



UNIVERSIDAD CARLOS III DE MADRID

ESCUELA POLITÉCNICA SUPERIOR

Ingeniería en Informática

Proyecto Fin de Carrera

**Análisis de la contratación en la
administración estadounidense para la
predicción de futuras licitaciones**

Autor: Francisco París Soriano

Director: Jose María Álvarez Rodríguez

Octubre 2016

Agradecimientos

La vida te lleva por caminos raros. Nunca pensé que pasarían tanto tiempo y tantas cosas entre el último día que pasé en la Universidad y el primero en que empezaría el proyecto final, así que lo justo en este momento sería agradecer a toda la gente que me ha ayudado y me ha enseñado algo en este camino: gracias a todos.

Sí quiero acordarme especialmente y escribir mi agradecimiento explícito para:

Sara, porque contigo la vida me mira con los labios pintados.

Lucía y Ana, mi pequeña y gran revolución.

Mis padres y hermana: raíz sobre todas las cosas.

Todos aquellos compañeros que estuvieron más cerca en la primera vuelta, especialmente a Carlos y David por mantenernos a flote en las peores tormentas.

Todos los que me acompañaron en los cursos de La Estación, donde todos aprendimos algo. Gracias Lucía y María José.

Mis ya excompañeros de Indra: Lore, Noe, Jesús, Zaloa, Ángel, Cristina, Miguel Ángel, Raquel, Charo, María, Rafa, Merche, Carmen...a todos los que pusieron un poquito para lograr un montón. Que la fuerza os acompañe.

Jose María, por su inestimable ayuda durante el desarrollo de este trabajo. No es habitual encontrar tan buenas ideas y ánimos tan cercanos, gracias de verdad.

Índice de contenidos

Resumen	9
1. Introducción.....	10
1.1 Objetivos.....	12
1.2 Estructura de la Memoria	13
2. Estado del Arte	15
2.1 Ficheros de Contratación	15
2.2 Big Data	16
2.3 Regresión múltiple lineal.....	21
2.4 Series Temporales y Modelos ARIMA	25
2.5 Redes Neuronales Artificiales. Modelo LSTM.	28
3. Análisis	33
3.1 Introducción.....	33
3.2 Definición de requisitos de usuario	33
3.3 Establecimiento de los requisitos software.....	36
3.4 Matriz trazabilidad requisitos usuario/SW	42
4. Diseño	43
4.1 Arquitectura de la solución. Diagrama de Clases.....	43
4.2 Entorno tecnológico.....	46
4.2.1 Apache Flink	46
4.2.2 Lenguaje R y R Studio	55
4.2.3 Keras y TensorFlow	56
4.2.4 Docker	57
4.2.5 Python.....	58
4.2.6 Java y Eclipse	59
5. Implementación	60
5.1 Introducción.....	60
5.2 Análisis inicial de los ficheros.....	60
5.3 Preselección de características.....	62
5.4 Estudio de las características preseleccionadas	65

5.5	División contratos.....	70
5.6	Análisis de la evolución de la contratación	71
5.7	Transformación de ficheros de nuevas oportunidades	75
5.8	Recuento	77
5.9	Completación de la serie temporal	78
5.10	Regresión lineal sobre múltiples predictores.....	80
5.11	Función auto.arima sobre series concretas	85
5.12	Redes LSTM.....	87
6.	Experimentación.....	92
6.1	Rendimiento en el tratamiento de datos	92
6.2	Experimentación agrupando por presupuesto y NAICS.....	93
6.3	Experimentación agrupando por presupuesto y reduciendo el NAICS	95
6.4	Comparativa de la aproximación de los modelos para un conjunto de series	97
6.5	Matriz de validación de requisitos SW	100
7.	Conclusiones.....	102
8.	Líneas Futuras	105
9.	Gestión del Proyecto.....	107
9.1	Planificación	107
9.2	Presupuesto.....	115
10.	Anexo I: Manual de Usuario	120
11.	Bibliografía.....	124

Índice de figuras

Figura 2.1 – Infografía “Internet en 1 minuto” (fuente: Excelacom Inc.).....	17
Figura 2.2 - Panorámica de herramientas y soluciones del ecosistema Big Data	18
Figura 2.3 - Esquema del funcionamiento del algoritmo MapReduce	19
Figura 2.4 - Logotipos de las empresas Cloudera, HortonWorks y de la solución Google Cloud Platform.....	21
Figura 2.5 - Ejemplo de uso de la función auto.arima en R	28
Figura 2.6 - Diagrama de elementos de una neurona artificial.....	29
Figura 2.7 - Esquema de una red neuronal artificial	30
Figura 2.8 - Esquema de una neurona LSTM.....	32
Figura 3.1 - Matriz trazabilidad Req. Usuario/SW.....	42
Figura 4.1 -Secuencia para el tratamiento de los ficheros.....	43
Figura 4.2 - Aplicación de modelos y selección del más adecuado	44
Figura 4.3 - Diagrama de Clases empleadas en la fase de implementación.....	45
Figura 4.4 - Logotipo de Apache Flink	46
Figura 4.5 - Arquitectura y componentes de Apache Flink.....	46
Figura 4.6 - Ejemplo programa en Apache Flink y cadena de transformaciones.....	48
Figura 4.7 - Vista de la interfaz web cliente con un mapa optimizado de ejecución	49
Figura 4.8 - Vista del cuadro de mando con el resultado de ejecución	49
Figura 4.9 - Vista principal del cuadro de mandos.....	50
Figura 4.10 - Logotipos de R y R Studio	55
Figura 4.11 - Logos de TensorFlow y Keras.....	56
Figura 4.12 - Logotipo de Docker	57
Figura 4.13 - Comparativa aplicación sobre máquina virtual y sobre Docker.....	58
Figura 4.14 - Logotipo de Python	58
Figura 4.15 - Logotipos de Java y Eclipse	59
Figura 5.1 - Extracto de un fichero de texto con información de los contratos	60
Figura 5.2 - Distribución de los contratos según su presupuesto	66
Figura 5.3 - Contratos con importes más altos y más bajos.	67
Figura 5.4 - Evolución del número de nuevos contratos	72
Figura 5.5 - Evolución del presupuesto destinado a nuevos contratos.....	73
Figura 5.6 - Evolución del número de contratos vs. modificaciones	74
Figura 5.7 - Evolución del presupuesto de nuevas oportunidades vs. modificaciones .	74

Figura 5.8 - Esquema explicativo de una línea del fichero transformado	77
Figura 5.9 - Histograma del número de series con observaciones distintas de 0	78
Figura 5.10 - Valor de los coeficientes de correlación entre las variables independientes a emplear en el modelo de regresión	81
Figura 5.11 - Parámetros del modelo de regresión múltiple lineal generado	82
Figura 5.12 - Resultados para la función <code>vif(Modelo_lm)</code>	82
Figura 5.13 - Gráfico que contrasta los valores ajustados frente a los residuos.....	83
Figura 5.14 - Gráfico Normal Q-Q.....	84
Figura 5.15 - Gráfico de Distancias de Cook	84
Figura 5.16 - Ejemplo aplicación función <code>auto.arima</code>	85
Figura 5.17 - Valores del modelo ARIMA de los criterios y de los coeficientes.....	86
Figura 5.18 - Valores de la predicción para 12 observaciones.....	86
Figura 5.19 - Ejemplo de implementación en Java del uso de <code>auto.arima</code> en R.....	87
Figura 5.20 - Esquema red de neuronas implementado en Keras/TensorFlow	88
Figura 5.21 - Implementación en Python empleando Keras	90
Figura 6.1 - Diagrama de Clases empleadas en la experimentación	94
Figura 6.2 - Regresión múltiple empleando sólo presupuesto y NAICS	95
Figura 6.3 - Regresión múltiple empleando sólo presupuesto y NAICS reducido	97
Figura 6.4 - Gráfico con el porcentaje de series que mejor cubre cada modelo.....	99
Figura 6.5 - Matriz Grado Cumplimiento Req. SW	100
Figura 9.1 - Tareas y diagrama de Gantt de la fase de planificación inicial	110
Figura 9.2 - Desglose de tareas de la planificación final.....	116
Figura 9.3 - Diagrama de Gantt con las tareas de planificación final	117

Índice de tablas

Tabla 2.1 - Número de registros incluidos en cada fichero	15
Tabla 3.1 - Requisito Usuario UR-C-01	34
Tabla 3.2 - Requisito Usuario UR-C-02	34
Tabla 3.3 - Requisito Usuario UR-C-03	34
Tabla 3.4 - Requisito Usuario UR-C-04	35
Tabla 3.5 - Requisito Usuario UR-C-05	35
Tabla 3.6 - Requisito Usuario UR-C-06	35
Tabla 3.7 - Requisito Usuario UR-C-07	35
Tabla 3.8 - Requisito Usuario UR-R-01	36
Tabla 3.9 - Requisito Usuario UR-R-02	36
Tabla 3.10 - Requisito Usuario UR-R-03	36
Tabla 3.11 - Requisito Software SR-F-01	37
Tabla 3.12 - Requisito Software SR-F-02	38
Tabla 3.13 - Requisito Software SR-F-03	38
Tabla 3.14 - Requisito Software SR-F-04	38
Tabla 3.15 - Requisito Software SR-F-05	38
Tabla 3.16 - Requisito Software SR-F-06	39
Tabla 3.17 - Requisito Software SR-F-07	39
Tabla 3.18 - Requisito Software SR-F-08	39
Tabla 3.19 - Requisito Software SR-F-09	39
Tabla 3.20 - Requisito Software SR-F-10	40
Tabla 3.21 - Requisito Software SR-O-01	40
Tabla 3.22 - Requisito Software SR-O-02	40
Tabla 3.23 - Requisito Software SR-O-03	40
Tabla 3.24 - Requisito Software SR-O-04	41
Tabla 3.25 - Requisito Software SR-R-01	41
Tabla 3.26 - Requisito Software SR-D-01	41
Tabla 3.27 - Requisito Software SR-D-02	42
Tabla 5.1 - Número de líneas por cada fichero de año fiscal	65
Tabla 5.2 – Valores discretos para los rangos de presupuesto	68
Tabla 5.3 - Número de nuevas contrataciones y de modificaciones	71
Tabla 5.4 - Serie de datos real y aproximada por ARIMA para el año 2015	86

Tabla 6.1 - Rendimientos para un sólo fichero.....	92
Tabla 6.2 - Rendimiento para la totalidad de los ficheros	92
Tabla 6.3 - Categorías generales para el NAICS.....	96
Tabla 6.4 - Series seleccionadas con observaciones mayores de 0	99
Tabla 6.5 - Resultados porcentuales sobre el total de series seleccionadas.....	99
Tabla 6.6 - Extracto del fichero generado en la clase de ComparativaModelo.....	101
Tabla 9.1 - Número de horas dedicadas a cada fase	115
Tabla 9.2 - Costes de recursos humanos.....	118
Tabla 9.3 - Costes de hardware.....	118
Tabla 9.4 - Costes de software y licencias.....	119
Tabla 9.5 - Resumen presupuesto total.....	119

Resumen

En el año fiscal de 2015, el Gobierno Federal de Estados Unidos ha manejado un presupuesto superior al billón de dólares y ha licitado durante el ejercicio más de dos millones y medio de contratos.

Existe una gran cantidad de datos sobre los contratos concedidos durante al menos el último quinquenio y es accesible a través de un portal gubernamental de forma pública y gratuita. La explotación de estos datos permite obtener una gran cantidad de información de cómo, cuándo, cuánto y dónde se distribuye el presupuesto, generando nuevas oportunidades que sepan aprovechar el valor de la información sintetizada.

Si se consigue predecir cuándo van a salir a concurso determinadas oportunidades de negocio se pueden anticipar los recursos humanos, económicos y materiales necesarios por parte de aquellas empresas que puedan estar interesadas en elaborar una oferta para el concurso público. Para ello, en este proyecto se van a procesar los ficheros de datos origen hasta conseguir la información que permita, mediante un modelo matemático o aprendizaje automático, describir cuándo se va a licitar el siguiente tipo de contrato atendiendo a unos parámetros definidos.

El reto de procesar la ingente cantidad de información disponible sobre los contratos es asumible gracias a utilidades de *big data* para el procesamiento de grandes volúmenes de datos. En este caso es Apache Flink el encargado de tratar los ficheros de contratos hasta transformarlos en series temporales. Esas secuencias de valores se van a abordar con tres modelos diferentes: regresión lineal, redes de neuronas y modelo ARIMA, evaluando el resultado obtenido con cada uno de ellos en función de su capacidad para representar la serie temporal de concesión de contratos.

1. Introducción

El ser humano ha demostrado un interés permanente por conocer el futuro con antelación a que éste ocurra desde tiempos remotos. En casi cualquier cultura y por casi cualquier método -por muy poco ortodoxos que hoy puedan parecernos-, la humanidad ha tratado de conocer en el presente lo que le espera en el mañana. A lo largo de la Historia es frecuente encontrar multitud de ejemplos en los que se trataba de discernir el futuro a partir del comportamiento de diferentes elementos naturales: efectuando sacrificios, observando el comportamiento de los astros, etc. A este tipo de adivinación, catalogada como artificial o inductiva, se le fue sumando y superponiendo una adivinación más deductiva, más intuitiva. Este segundo tipo de adivinación se acerca más a los términos de predicción o pronóstico que hoy manejamos: ya no se basa el pronóstico de lo que ocurrirá en la interpretación particular que un adivino hace de un determinado fenómeno natural sino en las observaciones de unos hechos ocurridos en el pasado que pueden ayudarnos a determinar que se vuelvan a dar las condiciones necesarias para que se repitan en el futuro de manera semejante.

El filósofo romano, Marco Tulio Cicerón, ya en el siglo I a.C. recogía en sus obras *De Divinatione* y *De oratore* [1] estas dos clasificaciones del arte adivinatorio que han llegado hasta nuestros días. Su famosa cita “*Historia est magistra vitae*” (“La Historia es maestra de la vida”) contiene la esencia de los fundamentos en que se va a basar este proyecto: hay que conocer y tener presente el pasado para determinar cómo incidirá éste en el futuro.

En la actualidad son muchos los campos donde se aplican diferentes técnicas de predicción para estimar un acontecimiento futuro, y el tratamiento de información histórica se ha convertido en una disciplina imprescindible en esos dominios donde poder evaluar los valores esperables en el futuro tiene un impacto directo en la cuenta de resultados. Se parte de la hipótesis de que los hechos pasados tienen un alcance significativo sobre los hechos observados en el presente y futuro, por lo que si se puede cuantificar esta influencia se podrá determinar el próximo hecho a ser observado.

Se emplea el término de ‘*serie temporal*’ para un conjunto de valores observados a lo largo de un período de tiempo y ordenados cronológicamente. Analizando el comportamiento de la serie, determinando sus ciclos, analogías, patrones de repetición, tendencia, estacionalidad, interrupciones y estudiando su posible relación con otros factores externos más allá de la simple repetición temporal de los valores, se trata de predecir el comportamiento futuro de la serie y el margen de error de esa predicción.

La secuencia de datos que conforma la serie temporal de estudio no siempre es fácil de observar y puede requerir un esfuerzo extraordinario para realizar su medición. En la actualidad se manejan cantidades ingentes de datos, más o menos estructurados, que suponen un reto en su tratamiento por su volumen y heterogeneidad. Pensemos, por ejemplo, que se quiere conocer el número de accesos por día de un usuario determinado

a una red social que cuenta con millones de usuarios activos durante los últimos tres meses y sólo se dispone para ello de los ficheros de log del servidor donde quedan registradas las entradas de todos los usuarios. Habría que almacenar todos esos ficheros, recopilarlos, conocer su estructura, tratarlos y ser capaz de procesar las millones de líneas que contienen para poder formar la serie temporal del ejemplo. Esto que hasta hace muy poco se convertía en un trabajo lento y oneroso si, por ejemplo, se cargaban esos ficheros en una base de datos relacional (RDBMS) y se trataban esos registros, se ha vuelto mucho más barato y sencillo con las nuevas herramientas de *Big Data* disponibles en el mercado.

Aplicaciones como Apache Hadoop [26], Apache Spark [29] o Apache Flink [25] permiten el tratamiento de datos masivos utilizando muchos menos recursos computacionales que las soluciones que existían hasta hace menos de un lustro. Empleando algunas de estas herramientas de forma eficiente se pueden determinar los valores de la serie temporal del ejemplo planteado mucho más rápido y con menos costes.

El auge del paradigma *Big Data* abre un nuevo camino de posibilidades en el tratamiento y procesamiento de información, permitiendo aplicar a grandes cantidades de datos distintos algoritmos de Aprendizaje Automático (del inglés, *Machine Learning*, *ML*) que sirvan para entender mejor la estructura, comportamiento e interrelaciones de los datos estudiados.

Una vez que se ha finalizado el tratamiento de los ficheros de partida del problema planteado y se tienen las mediciones de la serie temporal especificada es necesario encontrar un modelo que determine el comportamiento de ésta. El estadístico británico George E. P. Box, recoge en su libro *Empirical Model-Building and Response Surfaces* la famosa cita: “*Essentially, all models are wrong, but some are useful*”, (que se puede traducir como “*En esencia, todos los modelos están equivocados, pero algunos son útiles*”). La realidad es que cuando se trata de aplicar un modelo a una serie temporal ninguno resulta perfecto y no es capaz de ajustarse perfectamente a los datos, aunque sólo sea simplemente porque la serie de entrada contenga ruido en sus mediciones. Pero eso no significa que no puedan ser de utilidad y que permitan modelar la serie y elaborar predicciones sobre su comportamiento futuro con un grado de incertidumbre acotado.

Para el análisis y predicción de series temporales pueden emplearse multitud de técnicas; para la elaboración de este proyecto, de todas ellas se han escogido tres que se han empleado con éxito en bastantes casos anteriores: aplicación de regresión lineal, aplicación de modelos ARIMA y tratamiento de la serie con una red de neuronas artificial. Se estudiará el rendimiento de cada modelo escogido en el dominio del problema sobre predicción de contratación pública que se está tratando.

De un tiempo a esta parte, la opinión pública viene exigiendo medidas de mayor transparencia en la contratación realizada con dinero público que permita conocer mejor cómo y en qué se invierte el presupuesto de la administración, reclamando un mayor conocimiento del empleo de los recursos públicos, muchas veces con el ánimo de evitar

que se repitan los errores y abusos ocurridos en el pasado y presente más reciente. Un ejemplo de respuesta por parte de una administración pública -en este caso, estadounidense- a esa necesidad de información por parte de la ciudadanía es la Ley FFATA (*Federal Funding Accountability and Transparency Act*) de 2006, donde se exige que la información no confidencial de cada contrato firmado por la administración con una cuantía superior a los 25.000\$ se recoja en una web accesible a todos donde se pueda consultar y descargar esa información [20]. Los ficheros con toda la información contractual desde entonces están disponibles vía web y servirán de materia prima a este proyecto, tratándolos de manera automática con herramientas de *Big Data* y realizando predicciones a partir de esos datos.

1.1 Objetivos

Todo parte de un reto: conseguir predecir con la mayor exactitud posible los contratos que licitará la administración pública de EE. UU. el próximo año fiscal. Importa cuándo se va a firmar un contrato de una determinada tipología, atendiendo especialmente al presupuesto y al sector de actividad del servicio o producto requerido. Si se consigue saber con antelación y con una probabilidad de error aceptable cuándo se va a licitar un contrato, se puede informar a los posibles candidatos con antelación para que puedan ir trabajando en la elaboración de sus propuestas y tengan disponible los recursos humanos y materiales necesarios para presentar su candidatura en la forma y plazo requeridos. La materia prima de la que se parte son los ficheros de datos históricos [20], de dominio público, con la información de los contratos firmados desde el año 2000 hasta la actualidad.

Desde un punto de vista personal se buscaba que este proyecto supusiera un primer acercamiento al universo *Big Data*, por lo que -en la medida de lo posible- se buscaba priorizar el uso de herramientas de este contexto, de ahí la elección de Apache Flink como plataforma de procesamiento de los ficheros de datos históricos de contratos.

Para realizar el estudio de los datos ya procesados se ha optado por utilizar R [32] como lenguaje de referencia para el análisis estadístico de los mismos, por su uso generalizado entre los científicos de datos (del inglés, *Data Scientist*), lo que ha facilitado la proliferación de paquetes añadidos de uso libre con multitud de funcionalidades que facilitan la tarea de análisis. Para la implementación de un modelo de red de neuronas artificial se buscaba incluir en el proyecto alguna de las herramientas más novedosas, por lo que se decidió el empleo de TensorFlow [30] como herramienta de aprendizaje automático desarrollada por Google y de reciente liberación.

1.2 Estructura de la Memoria

En el capítulo 2 se presenta el estado del arte actual de las áreas de conocimiento relacionadas con este trabajo: ecosistema *Big Data*, series temporales utilizando modelos ARIMA y redes neuronales artificiales. Se incluye también la información básica necesaria para conocer el proceso de licitación y contratación de la administración pública estadounidense.

En el capítulo 3 se incluye el análisis de la solución a implementar, con la relación de requisitos de usuario y el conjunto de requisitos software definidos.

El capítulo 4 está dedicado al diseño de la solución, incluyendo el diagrama de la arquitectura de la aplicación, el diagrama de clases implementadas y la colección de aplicaciones y herramientas empleadas durante todo el proceso.

En el capítulo 5 se exponen las acciones realizadas en la implementación del proyecto para dar solución al problema planteado. Se explica la funcionalidad de las clases Java elaboradas en la fase de tratamiento y procesamiento de los ficheros de contratación y los modelos implementados sobre las series temporales de datos obtenidas, explicando los resultados de la aplicación de regresión lineal múltiple, modelos ARIMA y la red de neuronas artificial LSTM.

El capítulo 6 está destinado a la experimentación, donde se relatan las alternativas probadas para tratar de mejorar la eficiencia de los modelos y se realiza una comparativa entre ellos analizando los resultados obtenidos.

En el capítulo 7 se escriben las conclusiones alcanzadas y se reflexiona sobre el grado de cumplimiento del objetivo del proyecto.

El capítulo 8 recoge una serie de trabajos complementarios y líneas futuras a aplicar para mejorar el rendimiento de la solución propuesta, incrementar la calidad de aproximación de los modelos, y sugerencias para aplicar lo elaborado a otras fuentes de datos o tratar de capitalizarlo con fines comerciales.

El noveno capítulo se dedica a la gestión del proyecto incluyendo el desglose por tareas, su planificación al inicio del proyecto y el cumplimiento de las mismas al finalizar el período de trabajo. Se incluye también una pequeña memoria económica con la simulación del coste que hubiera tenido todo el proceso.

En el capítulo 10 se ha incluido un anexo con un manual de usuario donde se explican los ficheros de propiedades que recogen la configuración necesaria para la ejecución de las clases Java desarrolladas y los parámetros que se deben especificar para su correcta ejecución.

El capítulo 11 recoge las referencias a los libros, artículos y sitios web consultados que han servido como objeto de estudio para la elaboración de todo el proyecto y que son referenciados a lo largo de esta memoria de trabajo.

2. Estado del Arte

2.1 Ficheros de Contratación

Desde 2007 está disponible el sitio web UsaSpending.gov [20], donde se recogen ficheros con hasta 250 características de aquellos contratos concedidos por un valor superior a los 25.000\$ desde la administración pública estadounidense. Una parte importante de los contratos de menos de ese importe también están disponibles en esos mismos ficheros aunque no haya obligación legal de incluirlos.

Es posible descargarse los ficheros de los datos con los contratos de los años fiscales desde el 2000 hasta el 2016. El año fiscal transcurre desde el 1 de octubre del año natural hasta el 30 de septiembre del siguiente año. Así pues, los contratos del año fiscal 2008 comienzan en el mes de octubre de 2007 y finalizan en el mes de septiembre de 2008.

En la Tabla 2.1 se muestra el número de registros que contiene cada uno de los años fiscales en el momento de inicio de este proyecto (se incluyen modificaciones sobre los ficheros periódicamente, por lo que a la día de hoy algunos de los ficheros pueden presentar un número distinto en el número de sus registros). En total, los 16 ficheros de los años fiscales del 2000 al 2015 requieren más de 82 GBs para ser almacenados.

Fichero	Año	Totales (Con 1 línea de cabecera)
2000_All_Contracts_Full_20160115.csv	2000	594.598
2001_All_Contracts_Full_20160115.csv	2001	642.062
2002_All_Contracts_Full_20160115.csv	2002	830.599
2003_All_Contracts_Full_20160115.csv	2003	1.183.842
2004_All_Contracts_Full_20160115.csv	2004	2.002.014
2005_All_Contracts_Full_20160115.csv	2005	2.923.304
2006_All_Contracts_Full_20160115.csv	2006	3.797.162
2007_All_Contracts_Full_20160115.csv	2007	4.111.342
2008_All_Contracts_Full_20160115.csv	2008	4.504.816
2009_All_Contracts_Full_20160115.csv	2009	3.496.846
2010_All_Contracts_Full_20160115.csv	2010	3.538.782
2011_All_Contracts_Full_20160115.csv	2011	3.395.660
2012_All_Contracts_Full_20160115.csv	2012	3.116.220
2013_All_Contracts_Full_20160115.csv	2013	2.504.846
2014_All_Contracts_Full_20160115.csv	2014	2.512.863
2015_All_Contracts_Full_20160115.csv	2015	3.639.203

Tabla 2.1 - Número de registros incluidos en cada fichero

Gracias a esta disponibilidad de información sobre los procesos de contratación de la administración norteamericana y de la de otros países con una oferta de datos similar se han llevado a cabo iniciativas como Public Spending [21], o tesis como [22].

2.2 Big Data

En los últimos dos siglos, de la mano de una fuerte explosión demográfica acompañada de un desarrollo cultural y tecnológico, se han ido generando más y más datos susceptibles de ser almacenados. Se ha pasado de grandes bibliotecas repletas de libros a grandes centros de datos con miles de terabytes [2]. De hecho, se estima que este año, 2016, el tráfico de Internet supere por primera vez el zetabyte y alcance los 2 zetabytes en 2019 [4].

Considerando estas cifras no es de extrañar que cada cierto tiempo se alcance el límite de capacidad de almacenamiento y tratamiento de datos por los medios actuales del momento y se origine una nueva solución que sea capaz de satisfacer esas necesidades crecientes. La última ha sido el nacimiento del denominado ‘*Big Data*’ (podría traducirse como *datos masivos*, pero dada la generalización del término se mantiene su enunciado en inglés).

En 2001 se acuña por primera vez este término [3] y desde entonces se ha ido desarrollado su definición y sumando herramientas a su ecosistema. En esencia, *Big Data* hace referencia al hecho de almacenar, procesar y extraer conocimiento de datos a gran escala que no era posible realizar con sistemas de software y bases de datos convencionales. Primero fueron tres [5] y ahora son hasta siete [7] las V’s empleadas para definir la naturaleza de los datos masivos:

- **Volumen** (del inglés, *Volume*): hace referencia a la ingente cantidad de datos generados y almacenados para su posible tratamiento posterior desde diversas fuentes. Esta magnitud crece casi exponencialmente con los años y supone un reto la conservación y procesamiento de tal cantidad de información. El uso intensivo de Internet y de los dispositivos móviles, además del Internet de las Cosas (del inglés, *IoT, Internet of Things*) en auge, hace que la cantidad de datos intercambiados cada día a través de la red suponga un auténtico desafío en su tratamiento. En [6] se puede consultar una estadística de todo lo que ocurre en un solo segundo en Internet actualizado, y en la Figura 2.1 se incluye una curiosa infografía [14].
- **Velocidad** (del inglés, *Velocity*): los datos se generan de manera rápida, continua y en grandes cantidades. La respuesta para su almacenaje, procesado y análisis para sacarle el máximo rendimiento debe ser también lo más cercana al tiempo real posible.
- **Variedad** (del inglés, *Variety*): la heterogeneidad de las fuentes que originan los datos hacen que estos no siempre dispongan de una estructura clara que facilite su tratamiento.

- **Veracidad** (del inglés, *Veracity*): no se puede suponer que los datos sean siempre limpios y precisos. Puede haber datos que no cumplan la estructura que se les supone definida, que contengan duplicidades, etc., por lo que se hará necesario tratarlos antes de ser utilizados.

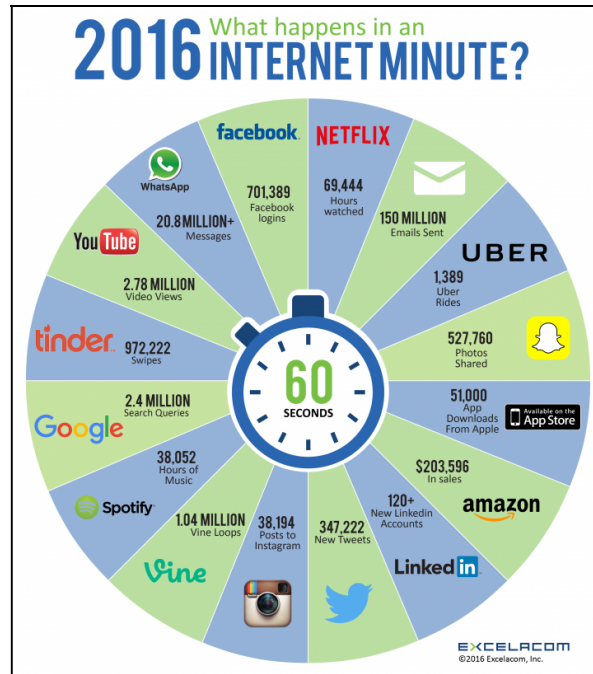


Figura 2.1 – Infografía “Internet en 1 minuto” (fuente: Excelacom Inc.)

- **Validación** (del inglés, *Validity*): hace referencia a la exactitud y precisión de los datos con respecto a la intención de uso.
- **Volatilidad** (del inglés, *Volatility*): los datos que existen y están disponibles ahora pueden no estarlo pasado un período de tiempo determinado, dado que pueden tener que eliminarse por requerimientos de espacio, legales, etc.
- **Valor** (del inglés, *Value*): todos los recursos que consumen la conservación y tratamiento de los datos obliga a que el aprovechamiento de éstos genere un retorno de la inversión de tal manera que sean capaces de generar información de valor para sus propietarios.

Algunos añaden a estas una más: **Visualización** (de diferentes términos en inglés: *Visualization*, *Visibility*), haciendo referencia a la capacidad de representación gráfica de la información generada tras el tratamiento de los datos.

Atendiendo a estas características no es de extrañar que las herramientas de *Big Data* surgidas no sólo se hayan centrado en el procesamiento y análisis de los datos sino que han abarcado todas las capas de una arquitectura tradicional en un sistema de *Business Intelligence* (BI, Inteligencia de Negocio): recolección de los datos, almacenamiento, el propio tratamiento y análisis de los datos y herramientas para su visualización y presentación.

En Figura 2.2 se incluye una panorámica actualizada de las herramientas que a día de hoy compiten y/o se complementan en el ecosistema *Big Data* [17].

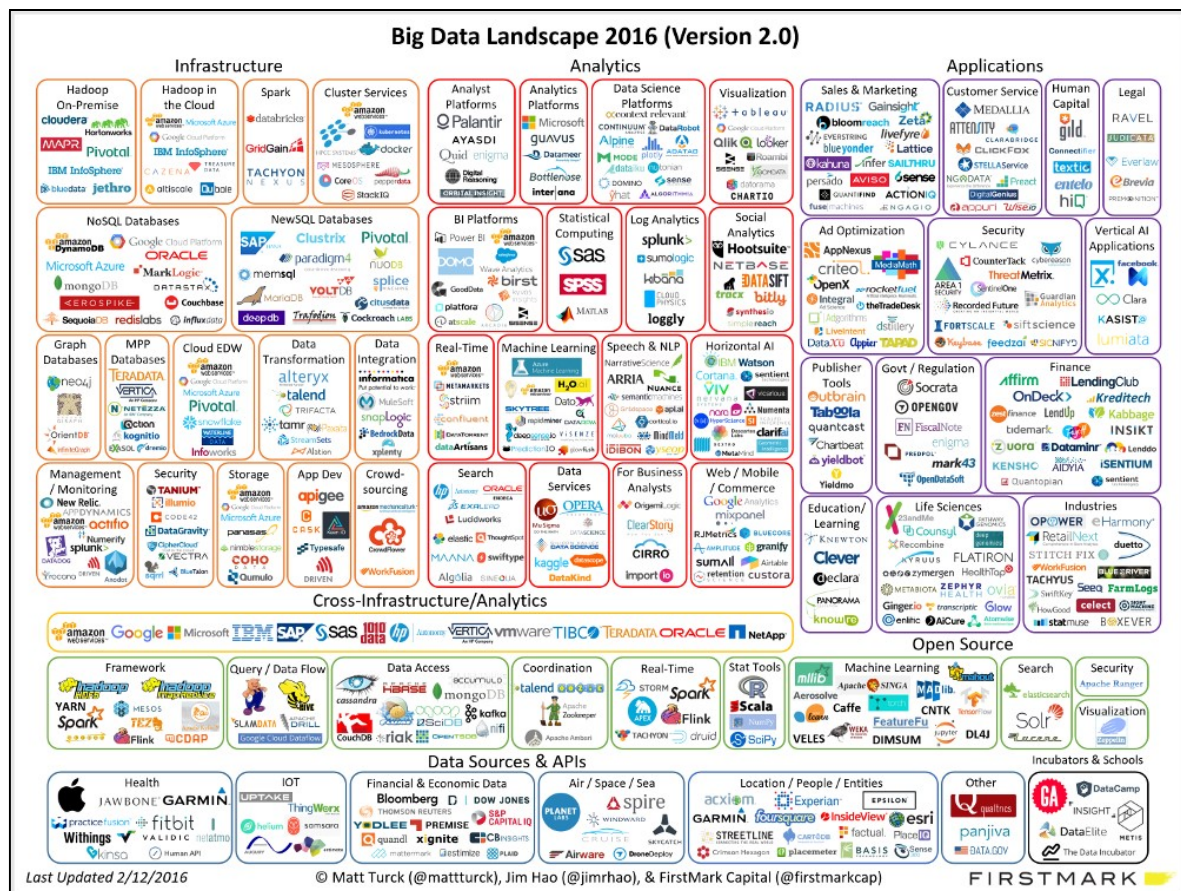


Figura 2.2 - Panorámica de herramientas y soluciones del ecosistema Big Data

Big Data ya no es considerada una tecnología emergente (de hecho, Gartner ya no la incluye en su último gráfico de 2015 [18] para *Hype Cycle for Emerging Technologies* aunque sí la incluía en el del 2014 [19]). Ha abandonado el *hype* y se están desinflando las sobredimensionadas expectativas que se tenían sobre ella para alcanzar después el grado de madurez y productividad de una tecnología consolidada.

□ Paradigma MapReduce

MapReduce es un modelo de programación concebido para satisfacer las necesidades de computación paralela sobre grandes ficheros de datos que tenía Google. A partir de entonces ha ido evolucionando a medida que se le incorporaban distintas variantes según la herramienta concreta que lo implementara, aunque sigue permaneciendo como concepto básico en todas ellas.

El nombre viene dado por la arquitectura del modelo, dividida principalmente en dos fases pensadas para ejecutarse en una infraestructura formada por múltiples nodos de manera distribuida, donde al menos uno de ellos desempeña un papel de responsable (*master*) y el resto desempeñan una tarea concreta (*workers*):

- **Map:** el nodo *master* es el encargado de particionar los datos de entrada en varios bloques a ser tratados en paralelo por los nodos *workers*. Cada partición es procesada en un nodo independientemente del resto, ejecutando sobre los datos una función de mapeo. El objetivo es generar a partir del bloque asignado a ese nodo un conjunto de pares clave-valor aplicando la función de mapeo determinada.

$$\text{Map}(\text{claveEntrada}_1, \text{valorEntrada}_1) \rightarrow \text{Lista}(\text{Clave}_2, \text{Valor}_2)$$

- **Reduce:** en la fase de reducción se reciben los pares generados en la fase de mapeo y se aplica una función de reducción sobre todos los valores asignados a una misma clave, obteniendo una nueva lista resultante. Las funciones reductoras más típicas son la agregación, sumatorio, máximo, mínimo, etc.

$$\text{Reduce}(\text{Clave}_2, \text{Lista}(\text{Valores}_2)) \rightarrow \text{Lista}(\text{Valor}_3)$$

Este paradigma de programación resulta poco flexible a la hora de intentar abordar diferentes tipos de problemas, ya que está limitado por la propia arquitectura del modelo; es por ello que algunas herramientas que se basan en él -como Apache Flink- implementen además otro tipo de transformaciones posibles sobre los datos de entrada: *joins*, iteraciones, etc.

En la Figura 2.3 puede verse un esquema sintetizado de la aplicación del algoritmo de *MapReduce* a un fichero de entrada donde se quiere calcular el número de apariciones de cada letra.

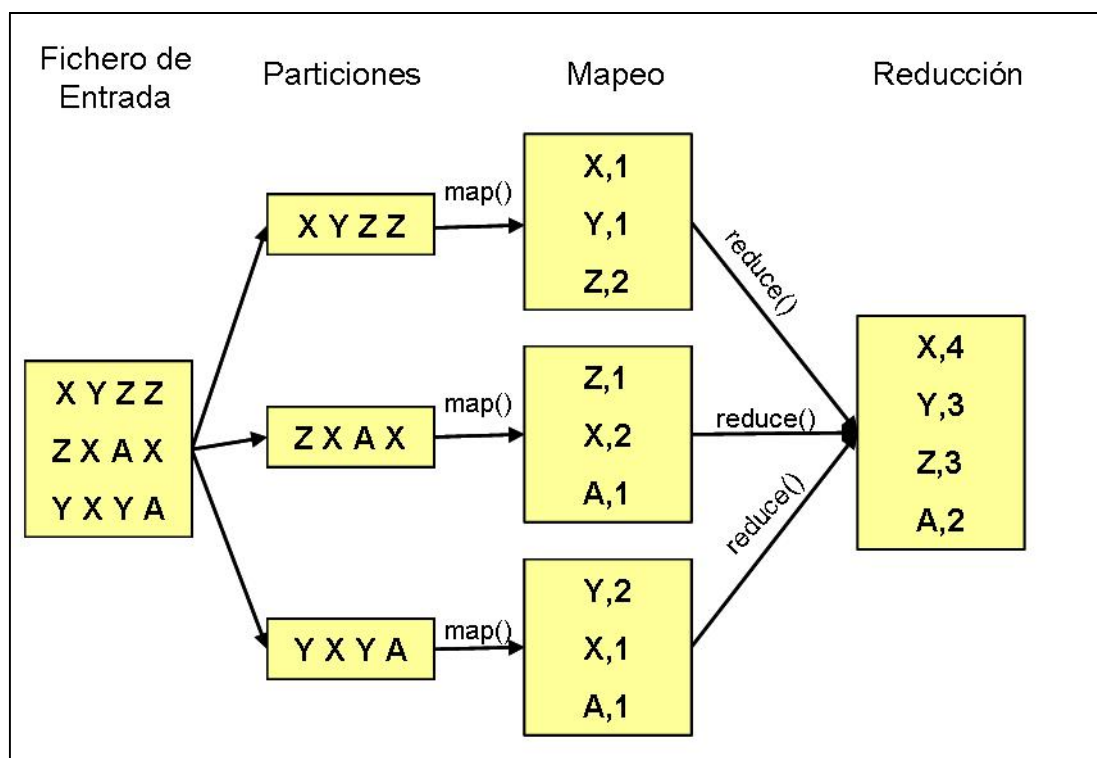


Figura 2.3 - Esquema del funcionamiento del algoritmo MapReduce

Apache Hadoop [26] y todas sus librerías relacionadas han sido durante los últimos años las herramientas más extendidas para lograr un exitoso tratamiento de grandes cantidades de datos. Utilizado en conjunto con soluciones de aprendizaje automático (por ejemplo, Apache Mahout [27]) y bases de datos NoSQL (por ejemplo, MongoDB [28]) han constituido una auténtica piedra de toque en la gestión de la información.

Apache Hadoop es una librería de código abierto (del inglés, *open source*) que facilita el procesamiento de grandes conjuntos de datos de manera distribuida entre clusters usando modelos sencillos de programar. Es escalable, puesto que permite su funcionamiento tanto en local como en grandes conjuntos de servidores, incluyendo un tratamiento de la gestión de fallos. Los módulos principales de Hadoop, en su versión 2.0, son:

- **Hadoop Common:** utilidades comunes necesarias para el resto de módulos.
- **Hadoop Distributed File System (HDFS):** provee un sistema de archivos distribuido accesible por las aplicaciones que se ejecuten en el clúster.
- **Hadoop YARN:** es el *framework* que se utiliza para la planificación de trabajos y la gestión de recursos del clúster.
- **Hadoop MapReduce:** permite el procesamiento en paralelo de grandes *datasets* (conjuntos de datos).

Más recientemente se ha popularizado el uso de Apache Spark [29] como alternativa más eficaz a Hadoop MapReduce. Promete ser hasta 100 veces más rápido que Hadoop cuando la ejecución es en memoria y hasta 10 veces más rápido cuando se realiza teniendo que volcar a disco.

Spark nace en 2009 en el AMPLab de la Universidad de Berkeley [8] y es donando a la fundación Apache en 2013. Hadoop MapReduce estaba concebido para un flujo de datos lineal: se leen de disco los datos, se aplica la función *map* sobre los datos, después se aplica la función *reduce* sobre los resultados y se vuelca de nuevo a disco. Spark permitía mantener en memoria el conjunto de datos y aplicarle los distintos *map-reduce* tantas veces como se necesitara sin necesidad de volcar a disco tras cada transformación, ahorrando una gran cantidad de tiempo que *MapReduce* empleaba en acceso a disco para lectura-escritura.

Tanto Apache Hadoop como Apache Spark son ampliamente utilizados por multitud de empresas como sus soluciones para tratamiento de grandes volúmenes de datos. Muchas de ellas, en lugar de hacer un uso directo de este software, emplean plataformas comerciales que implementan soluciones basadas en Hadoop/Spark, como por ejemplo: Cloudera [36], que ofrece la distribución de CDH (*Cloudera Distribution Hadoop*) que contiene el *core*, los principales elementos de Hadoop que proporcionan un procesamiento de datos fiable y escalable (MapReduce, HDFS), así como otros componentes específicos para empresas que aportan seguridad, alta disponibilidad e

integración con hardware y otras aplicaciones software. Otra alternativa muy popular es Hortonworks [38], que ofrece su plataforma *Hortonworks Data Platform* (HDP), similar a la anterior. Ambas soportan también la instalación de Apache Flink. Google también ofrece una distribución en la nube, *Google Cloud Platform*, con una potente API para gestionar -entre otras muchas cosas- soluciones que impliquen *Big Data*, permitiendo capturar, procesar, almacenar y analizar los datos dentro de la misma plataforma, que además ofrece también soluciones de aprendizaje automático integradas.



Figura 2.4 - Logotipos de las empresas Cloudera, HortonWorks y de la solución Google Cloud Platform

2.3 Regresión múltiple lineal

La aplicación de regresión múltiple [61][64], está motivada por la existencia de varias variables que están conectadas entre sí para formar una respuesta; a cada una de esas variables independientes se le denomina predictor.

Mientras que en la regresión lineal simple es una única variable independiente la que influye en la salida, en la regresión múltiple se emplea más de una variable independiente como predictor de la salida, atendiendo a la fórmula:

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k$$

Donde \hat{y} es la salida estimada en función de la combinación del término independiente β_0 y x_i representa a cada variable independiente hasta un total de k predictores. Lo que se calcula es cada uno de los términos β_i , minimizando la suma de los errores cuadráticos (*SSE*, *Sum of Squares of Error*).

La fórmula del *SSE* es la siguiente, siendo y_i el valor real e \hat{y}_i el valor estimado:

$$SSE = \sum_{i=1}^k (y_i - \hat{y}_i)^2$$

Para el cálculo de los coeficientes β_i de la fórmula de la regresión se suelen emplear soluciones computacionales: por ejemplo, el cálculo se realiza automáticamente con la función *lm* (*linear model*) del paquete *stats* en el entorno de R, que se utilizará en la fase de implementación.

Uno de los parámetros que evalúa la bondad del modelo obtenido y de cómo este se ajusta a los datos es el Coeficiente de Determinación, R^2 , entendido como la proporción de la varianza de los datos que es explicada por el modelo.

Para el cálculo de R^2 , además del SSE , es útil conocer el SST y el SSR :

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}, \text{ siendo } SSR = \sum (\hat{y} - \bar{y})^2 \text{ y } SST = \sum (y - \bar{y})^2$$

Siendo SST (Suma de Cuadrados de Tratamiento, del inglés *Sum of Squares Total*), entendida como la suma del cuadrado de la diferencia de los valores reales respecto de la media, y SSR (Suma de los cuadrados de la regresión, *Sum of Squares Regression*) igual que la anterior pero aplicado a los valores estimados con el modelo.

Esta forma de calcular R^2 es aceptable cuando se trata de un único predictor, pero cuando se emplean varios, como es este caso de la regresión múltiple, es necesario ajustarlo para que se tenga en consideración la varianza explicada por la totalidad de los predictores empleados.

La fórmula del R^2_{adj} es la siguiente:

$$R^2_{adj} = 1 - (1 - R^2) \times \frac{n - 1}{n - k - 1}$$

Siendo n el número de muestras utilizadas en el cálculo del modelo y k el número de predictores (variables independientes) empleadas en el cálculo del mismo.

Este indicador R^2_{adj} se utilizará en la fase de implementación de construcción del modelo de regresión múltiple como indicador de idoneidad del modelo (R^2_{adj} toma valores entre 0 y 1: cuanto más cerca está del 1, más varianza de los datos se explica según el modelo obtenido).

Es importante encontrar el modelo que mejor se ajuste a los datos porque así, en principio, si el conjunto de observaciones posterior a predecir no está muy alejado de la realidad de los datos con los que se ha construido el modelo, se va a obtener la mejor estimación. El mejor modelo no siempre es el más complicado y no tiene por qué incluir todas las variables independientes como predictores. Atendiendo a la evolución del coeficiente de determinación ajustado se pueden emplear dos formas distintas de encontrar el mejor modelo:

- **Inclusión de predictores** (del inglés, *Forward Selection*): se comienza añadiendo una a una las variables independientes como predictores al modelo y se recalcula el coeficiente: si éste es mejor entonces la variable se incluye en el modelo; si no mejora, se descarta y se prueba con otra variable.

- **Eliminación de predictores** (del inglés, *Backward Elimination*): se comienza con un modelo completo donde están incluidas todas las variables independientes como predictores. Se van eliminando una a una: si el coeficiente mejora, esa variable se descarta; si no mejora, se mantiene esa variable y se prueba a eliminar la siguiente y recalculan el modelo.

No siempre se llega al mismo modelo resultado empleando ambas técnicas, así que si se busca encontrar el mejor de todos por coeficiente de determinación, hay que probar con las dos y luego comparar los valores obtenidos.

Para que un modelo de regresión múltiple lineal se pueda decir que constituye una forma adecuada de representación de una muestra de datos, además del coeficiente de determinación, debe cumplir en grado suficiente las siguientes características:

- Los residuos –diferencias entre el valor real y el ajustado- del modelo deben seguir una distribución normal. No deben presentarse irregularidades o valores extremos.
- La variabilidad de los residuos debe ser aproximadamente constante. Esto hace referencia a la homocedasticidad, buscando que todos los valores estén igual de bien aproximados y que no haya grandes diferencias.
- Los residuos deben ser independientes unos de otros.
- Cada variable independiente muestra una relación lineal con la variable de salida.

Para comprobar sobre un modelo obtenido cómo se cumplen estas condiciones, se emplean algunos gráficos de utilidad:

- Curva normal de errores (*Normal Probability Plot*), para comprobar que no se aprecian irregularidades.
- Un gráfico que representa los valores ajustados frente a los residuos, para comprobar que no existen desviaciones.
- Diagrama de Distancia de Cook: donde se aprecia cómo influyen las observaciones en la generación del modelo resultante y se puede detectar aquellas que tienen una influencia excesiva en el resultado. Para este cálculo se suele calcular un valor de corte, que típicamente es de 4 dividido entre el número de grados de libertad (el número de grados de libertad viene dado por el número de muestras empleadas menos el número de variables utilizadas como predictores). Al generar el gráfico se representa en qué grado cada observación supera ese umbral y se pueden identificar valores extremos.

Si con estas utilidades se observa que las premisas que se han asumido no se cumplen, lo más probable es que los datos no se pueden representar adecuadamente con un modelo de regresión lineal múltiple y se debe buscar un modelo alternativo más adecuado.

□ Elaboración de un modelo de regresión múltiple en R

Antes de construir el modelo resulta útil estudiar el grado de correlación entre las variables a utilizar. Para ello se puede utilizar la función *cor* [63] de la librería *stats* de R: indicando como parámetro el *dataframe* y las variables a correlacionar es suficiente. Se puede indicar también qué coeficiente de correlación calcular (el de Pearson se utiliza por defecto). Hay otras alternativas extendidas como es el uso de la función *corr.test* del paquete *psych* que además ofrece más información complementaria.

Para elaborar el modelo de regresión la función a emplear es *lm* [62]. Se indica la variable independiente y las variables dependientes de las que va a depender, y se incluye la muestra de datos a partir de la que elaborar el modelo. La fórmula genérica de invocación sería esta:

$$\text{Modelo} = \text{lm}(y \sim x_1 + x_2 + \dots x_k, \text{data} = \text{dataframe})$$

Siendo *y* la variable dependiente a calcular en función de las variables independientes x_i .

Se pueden obtener los detalles del modelo resultante aplicando la función *summary* al modelo obtenido en la anterior ejecución, y ver así los coeficientes aplicados a cada variable independiente del modelo. Con *summary(Modelo)\$adj.r.squared* se devuelve el valor de coeficiente de determinación ajustado R^2_{adj} . Si además se quiere obtener el intervalo de confianza al 95% para cada coeficiente, se puede emplear para ello la función *confint(Modelo)*.

Existe además una función *vif(Modelo)* en la librería *car* que permite obtener el *Variance Inflation Factor (VIF)* de las variables predictores que participan en un modelo. Si una variable obtiene un valor de 5 o superior en este test se puede eliminar por considerarse redundante respecto al resto de variables del modelo.

En R se pueden dibujar los gráficos que se necesitan para comprobar las premisas que se enunciaban a la hora de proponer un modelo de regresión múltiple lineal:

- Con *plot(Modelo, which=1)* se puede visualizar gráficamente la relación entre los residuos y los valores ajustados.
- Para el diagrama de distancias de Cook donde se muestra el grado de influencia de cada observación, primero hay que construir el umbral, que típicamente es el resultado de dividir 4 entre los grados de libertad del modelo (*umbral <-*

`4/(Modelo$df)`) y generar el gráfico con la instrucción `plot(Modelo, which=4, cook.levels=umbral)`.

2.4 Series Temporales y Modelos ARIMA

[65][66] El estudio de una serie temporal tiene por objeto analizar la evolución de una variable a través del tiempo. La diferencia esencial entre las series temporales y los análisis no temporales (estadística descriptiva, regresión) es que en los análisis previos no importaba el orden en que estaban tomadas las observaciones y éste se podía variar sin problemas. En series temporales el orden es muy importante y variarlo supone cambiar la información contenida en la serie.

A la hora de identificar la serie temporal es necesario conocer:

- La **periodicidad** con la que están tomados los datos. Por ejemplo, se tiene una periodicidad mensual si se toma una observación de la serie cada mes.
- La **tendencia** que muestra la serie: si a largo plazo crece o decrece se dirá que posee tendencia positiva o negativa, respectivamente.
- La **variabilidad** de la serie: si es homocedástica (presenta una variabilidad más o menos homogénea a lo largo del tiempo) o heterocedástica, si presenta cambios reseñables en su variabilidad a lo largo del tiempo.
- **Estacionalidad**: si la serie presenta ciclos estacionales que se repiten en el tiempo.

Para realizar el estudio de una serie temporal se analiza el gráfico de la serie, donde quizás se podrán apreciar a simple vista algunas de las características de ésta enunciadas anteriormente.

Se calculan también la **Función de Autocorrelación Simple** (FAS) y la **Función de Autocorrelación simple Parcial** (FAP). La primera proporciona la estructura de dependencia lineal de la serie (cómo una observación influye en las siguientes) y con la segunda se puede ver cómo una observación proporciona la relación directa que existe entre observaciones separadas por k retardos entre ellas.

Un tipo especial de serie estacionaria es la serie denominada *ruido blanco*. Un ruido blanco es una serie estacionaria tal que ninguna observación influye sobre las siguientes. Lo que se busca en una serie es que no sea estacionaria, por lo que se elimina la tendencia y los ciclos estacionales restando diferencias a la serie.

Los procesos autorregresivos forman una familia de procesos tales que una observación depende de las observaciones anteriores. Se denominan procesos AR y se caracterizan por su orden. Tienen una característica común: todos ellos tienen memoria larga. La memoria larga, que se aprecia fácilmente al dibujar la FAS, se traduce en que

la serie tarda bastante tiempo en absorber una alteración extrema. Para aquellas series en las que se requiere que ese impacto externo se registre de forma inmediata se introduce una nueva familia de modelos para poder representarlos. La familia de procesos matemáticos que representa los procesos de memoria corta se denomina genéricamente procesos de media móvil, MA.

En la práctica se observa que las series no son puras AR o MA, sino que presentan parte AR y parte MA. Los procesos ARMA(p, q) son combinación de estructuras autorregresivas y de media móvil que tienen una parte AR(p) y una parte MA(q). Su FAS y su FAP serán combinación de ambos procesos. Los modelos ARMA únicamente sirven para ajustar series estacionarias.

Cuando se tiene una serie no estacionaria es necesario transformarla en una serie estacionaria para poder modelarla. Este tipo de series se modelan con ARIMA(p, d, q): donde p representa el orden de la parte autorregresiva de la serie estacionaria, q representa el orden de la media móvil de la serie estacionaria y d representa el número de diferencias que ha habido que tomar para que la serie que, inicialmente no era estacionaria, sea finalmente estacionaria.

El ajuste de un modelo ARIMA(p, d, q) implica cumplir satisfactoriamente la etapa de identificación del modelo. Si una serie está bien identificada, y se le aplica el modelado con ARIMA, los residuos de la serie deben carecer completamente de estructura. Como se ha dicho anteriormente, una serie sin estructura de dependencia es precisamente el ruido blanco.

La metodología Box-Jenkins sistematiza el cálculo de los parámetros del modelo tratando de obtener el mejor ajuste y buscando así que los pronósticos futuros sean más precisos. La fórmula matemática que representa un modelo ARIMA(p, d, q) es:

$$Y_t = -(\Delta^d Y_t - Y_t) + \phi_0 + \sum_{i=1}^p \phi_i \Delta^d Y_{t-i} - \sum_{i=1}^q \theta_i \varepsilon_{t-i} + \varepsilon_t$$

Considerando que la diferencia de observaciones es $\Delta Y_t = Y_t - Y_{t-1}$, que d corresponde a las diferencias a aplicar a la serie para convertirla en estacionaria, que los términos de parte autorregresiva (AR) son ϕ_i (desde 1 a p), que los parámetros de la parte de medias móviles (AM) son θ_i (desde 1 hasta q), ϕ_0 es una constante y ε_t es el término de error.

□ Modelos ARIMA en R

En R, el cálculo de un modelo ARIMA sobre una serie temporal univariada puede hacerse utilizando la función *arima* del paquete *stats*. A esta función, además de pasarle la serie de observaciones, hay que indicarle dos parámetros esenciales del modelo:

```
arima(x, order = c(0L, 0L, 0L), seasonal = list(order = c(0L, 0L, 0L),  
period = NA))
```

El parámetro *order* hace referencia a la parte no estacional del modelo y hay que indicar los valores para (p, d, q) . Lo mismo ocurre con *seasonal*, que es la parte estacional del modelo incluyendo el periodo de frecuencia.

Para aplicar esta función y obtener el mejor modelo, por tanto, es necesario haber identificado correctamente la serie y los parámetros que definan su parte no estacional y estacional. Existe una variante de esta función si no se dispone de estos parámetros y se quiere que se autocalcule el mejor modelo para la serie. Para ello se puede utilizar la función *auto.arima* del paquete *forecast* [68].

Esta función *auto.arima* [67] consiste en evaluar los modelos ARIMA susceptibles de ajustarse a la serie de observaciones pasada por parámetro y determinar en base a un criterio de selección cuál presenta mejor ajuste.

Los criterios de selección que permite son: [64] AIC (Criterio de Información de Aikake; del inglés *Akaike Information Criterion*) que es una medida de la calidad relativa de un modelo estadístico para un conjunto dado de datos. Como tal, el AIC proporciona un medio para la selección del modelo. AIC ofrece una solución de compromiso entre la bondad de ajuste del modelo y su complejidad. Basándose en la entropía de información ofrece una estimación relativa de la información perdida cuando se utiliza un modelo determinado para representar el proceso que genera los datos. AIC no valida o rechaza una hipótesis nula acerca del modelo, por lo que no constituye una medida acerca de la calidad del modelo en un sentido absoluto. Si todos los modelos candidatos encajan mal, AIC no dará ningún aviso de ello.

AICc (AIC corregido; del inglés, *AIC with a correction*) es similar al anterior pero incluyendo un factor de corrección para muestras de datos finitas basándose en el número de datos de la muestra y en el número de parámetros empleados en el modelo. BIC (Criterio de Información Bayesiana; del inglés, *Bayesian Information Criterion*) es también un criterio de selección de modelos que penaliza en mayor medida que el AIC el número de parámetros empleados en el modelo.

A la hora de ejecutar la función *auto.arima* se le pueden indicar los criterios a seguir, utilizando el parámetro *ic* de la función informado con cualquiera de los siguientes valores: *c("aicc", "aic", "bic")*. Para que en la salida se muestren todos los modelos ARIMA evaluados hay que indicar el parámetro *trace* a *TRUE*.

Una vez seleccionado el modelo, se utiliza la función *forecast* para predecir N número de observaciones siguientes a las facilitadas en la muestra de datos según el modelo calculado.

Un ejemplo de construcción del modelo y de la generación de las 3 siguientes observaciones se puede observar en la Figura 2.5.

```

Console ~/
> library(forecast)
> serieDatos <- c(10,25,30,40,38,27)
> modeloARIMA <- auto.arima(serieDatos, trace=T)

ARIMA(2,0,2) with non-zero mean : Inf
ARIMA(0,0,0) with non-zero mean : 52.46865
ARIMA(1,0,0) with non-zero mean : 61.27563
ARIMA(0,0,1) with non-zero mean : 60.94653
ARIMA(0,0,0) with zero mean      : 60.83941
ARIMA(1,0,1) with non-zero mean : 90.64216

Best model: ARIMA(0,0,0) with non-zero mean

> observacionesPrediccion <- forecast(modeloARIMA,h=3)
> observacionesPrediccion
  Point Forecast    Lo 80    Hi 80    Lo 95    Hi 95
7      28.33333 14.51474 42.15192  7.199625 49.46704
8      28.33333 14.51474 42.15192  7.199625 49.46704
9      28.33333 14.51474 42.15192  7.199625 49.46704
>

```

Figura 2.5 - Ejemplo de uso de la función auto.arima en R

2.5 Redes Neuronales Artificiales. Modelo LSTM.

En 1943, W. McCulloch y W. Pitts [54] presentan un modelo simplificado de neuronas donde estas se representaban como redes biológicas dentro de un modelo conceptual de componentes para circuitos que podían desempeñar tareas computacionales.

Una neurona, en términos biológicos, es la célula fundamental del sistema nervioso encargada de transmitir los impulsos nerviosos. Cada una consta de tres partes principales: el cuerpo de la célula (soma), donde se encuentra el núcleo de ésta y se realiza la síntesis de proteínas. Las dendritas, que son unas ramificaciones por las que la neurona es capaz de recibir los impulsos de otras neuronas con las que está conectada, sirviendo de entrada a la célula. Si la combinación de señales recibida por las dendritas es suficiente, se genera un impulso que se transmite por el axón de la neurona hasta las ramificaciones que constituyen la salida de la célula. Estas ramificaciones nerviosas de una neurona se unen con las dendritas de otras formando unas uniones denominadas sinapsis.

El equivalente artificial a una neurona biológica dentro de una red neuronal artificial (*ANN*, *Artificial Neuronal Network*) es el elemento procesador que se muestra en la Figura 2.6.

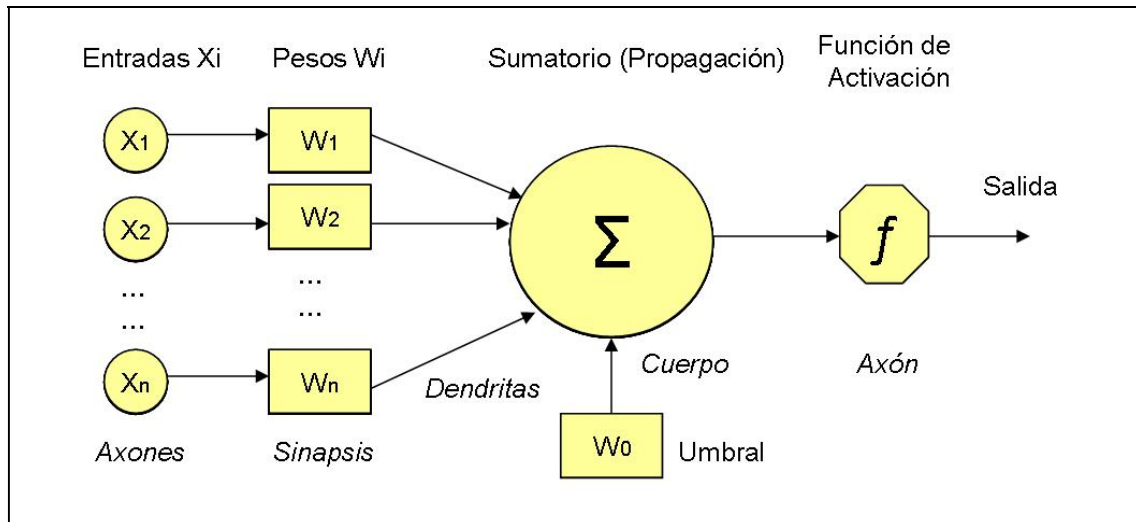


Figura 2.6 - Diagrama de elementos de una neurona artificial

Las sinapsis neuronales se modelan artificialmente como un conjunto de pesos aplicados sobre la señal de entrada a esa neurona. Un peso constituye un factor numérico: si este es positivo, constituye una conexión diseñada para excitar esa entrada de la neurona, mientras que si es negativo representa una inhibición.

A cada entrada se le aplica su factor de peso correspondiente y se agrega el total generado por todas ellas. A esta suma lineal de todas las entradas multiplicadas por su factor de peso se le aplica una función de activación para determinar la salida de la neurona en función de un umbral determinado, lo que decide la amplitud de la señal de salida que genera la neurona al resto que estén conectadas a ella.

□ Red de Neuronas artificiales

La potencia y capacidad de estos modelos artificiales –de manera análoga al biológico- reside en la capacidad de agrupación y conexión de estos elementos entre sí formando lo que se denominan redes. Las redes de neuronas se componen de un conjunto de elementos individuales como el descrito anteriormente conectado con otros similares formando capas, y estas capas a su vez se conectan con otras formando un entramado de elementos que forman la red. Se denomina *capa de entrada* a la constituida por las neuronas que reciben los valores de entrada a la red y *capa de salida* al conjunto de neuronas que generan los valores resultantes del procesamiento por parte de la red; las capas ocultas entre la capa de entrada y la de salida se denominan *capas intermedias u ocultas*. El dimensionamiento en número de neuronas de cada capa y la manera de interconectarlas determina el comportamiento de la red.

En la Figura 2.7 se incluye un esquema de la arquitectura de una red neuronal artificial (RNA) multicapa básica (perceptrón simple).

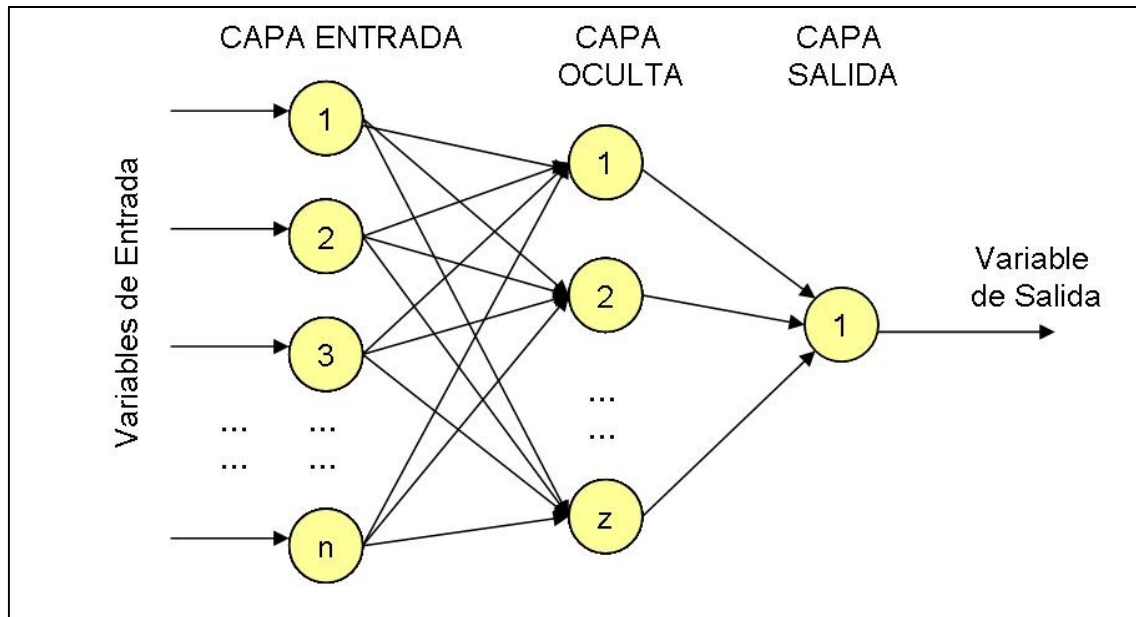


Figura 2.7 - Esquema de una red neuronal artificial

□ Ventajas y desventajas del empleo de una RNA

Los principales beneficios que presenta el empleo de redes artificiales son:

- **Aprendizaje adaptativo:** habilidad de aprender durante la fase de entrenamiento, ya sea por aprendizaje supervisado o no, ofreciendo un modelo que se ajuste a los datos utilizados. Se dice que es adaptativa porque no necesita ser reprogramada por un tercero: ella misma se va adecuando a la nueva realidad de los datos encontrados, reorganizando sus parámetros de forma autónoma. Durante el proceso de aprendizaje de una red neuronal se pueden aplicar dos principales diferentes modelos de aprendizaje:

- **Aprendizaje Supervisado:** el entrenamiento se realiza con un conjunto de datos de entrada etiquetados junto a la salida esperada. Se observa la salida de la red para los datos de entrada y se evalúa el error obtenido respecto a la salida correcta, ajustando los pesos y el umbral de activación según corresponda.
- **Aprendizaje No Supervisado:** no se precisa de un conjunto de valores de entradas etiquetados junto a un valor de salida porque no se conoce el valor de salida esperado. Se busca que la estructura determine patrones de comportamiento dentro de los datos y sea capaz de determinar las características más representativas.
- **Tiempo real:** los cálculos que realiza una red artificial ya entrenada pueden paralelizarse, por lo que de utilizarse un dispositivo que permita esta posibilidad pueden efectuarse en paralelo y ofrecer una salida en tiempo real.

- **Tolerancia a fallos:** la información dentro de la red se codifica de manera redundante así que de verse ésta destruida parcialmente es capaz de continuar ofreciendo una respuesta aceptable aunque se vea afectada por cierta degradación.
- **Aprendizaje profundo** (del inglés, *deep learning*): es un conjunto de algoritmos en aprendizaje automático que intentan modelar abstracciones de alto nivel en datos usando arquitecturas compuestas de transformaciones no-lineales múltiples. Las RNA son aproximadores funcionales universales, se ha demostrado que pueden aproximar cualquier función continua a cualquier precisión.
- **Flexibilidad:** una RNA puede desempeñar tareas que un programa lineal no y es capaz de gestionar el ruido en los datos de entrada minimizando el impacto en la salida frente a la información valiosa.

El principal inconveniente del empleo de un RNA es que requiere de una fase de entrenamiento que se vuelve bastante elevado cuanto mayor es el número de elementos de la red y más compleja es su estructura. Además, se necesita de un conjunto de datos de entrenamiento representativo y suficiente para lograr que la red sea capaz de generalizar sin sobreajustarse (del inglés, *overfitting*) a los datos concretos.

□ Redes Prealimentadas vs. Redes Recurrentes. Modelo LSTM.

Las redes prealimentadas (en inglés, *Feed-forward neural networks*) se caracterizan porque entre las neuronas todas las conexiones que se establecen son lineales y no se forman ciclos, de tal manera que una neurona sólo está conectada con otras de su misma capa o de capas posteriores, no existen conexiones hacia atrás y la información dentro de la estructura sólo fluye hacia adelante. Ejemplos típicos de redes de este tipo son Adaline y el Perceptrón.

Las redes recurrentes (del inglés, *Recurrent neural Networks, RNN*) sí pueden contener conexiones que acaben formando bucles entre neuronas de diferentes capas, sin importar que sean de capas anteriores. La principal ventaja que presentan frente al modelo previo es su capacidad de memoria en tanto que permiten a la información persistir a lo largo del flujo entre las distintas conexiones, de tal manera que entradas relevantes de pasos anteriores pueden llegar a evaluarse en un paso muy posterior. Por ello está siendo muy popular su empleo en tareas de reconocimiento del lenguaje y traducción donde una palabra depende del contexto y de varias palabras anteriormente encontradas, no sólo de la inmediatamente anterior.

Dentro de las redes recurrentes, las redes LSTM (redes con memoria a largo y corto plazo; del inglés, *Long Short-Term Memory*) [55] constituyen un tipo de red neuronal artificial recurrente especializado que se adapta muy bien en dominios donde a partir de los datos de entrada se deben clasificar, procesar y precedir series temporales que pueden presentar grandes intervalos de tiempo entre eventos representativos.

Cada neurona de una red LSTM tiene la capacidad de recodar un valor aparecido hace un tiempo indeterminado anterior en la serie; esto se consigue implementando en la neurona entradas que determinan cuándo la señal recibida es importante y debe ser

tenido en cuenta su recuerdo, cuándo se debe propagar este recuerdo o cuándo se debe descartar porque ya no es significativo.

Esta tipología de red surge para solucionar el problema que sufren las RNN en la fase de retropropagación (del inglés, *backpropagation*) por el llamado efecto de gradiente evanescente, cuando se van a ajustar los pesos en función del error obtenido. Durante esta fase la señal de gradiente puede ser multiplicada multitud de veces por la matriz de pesos asociadas a las conexiones lo que impacta en el aprendizaje, pudiendo darse el caso de que este se vuelva lento, deje de funcionar por completo o empiece a haber divergencias.

La estructura básica de cada neurona LSTM se muestra en la Figura 2.8.

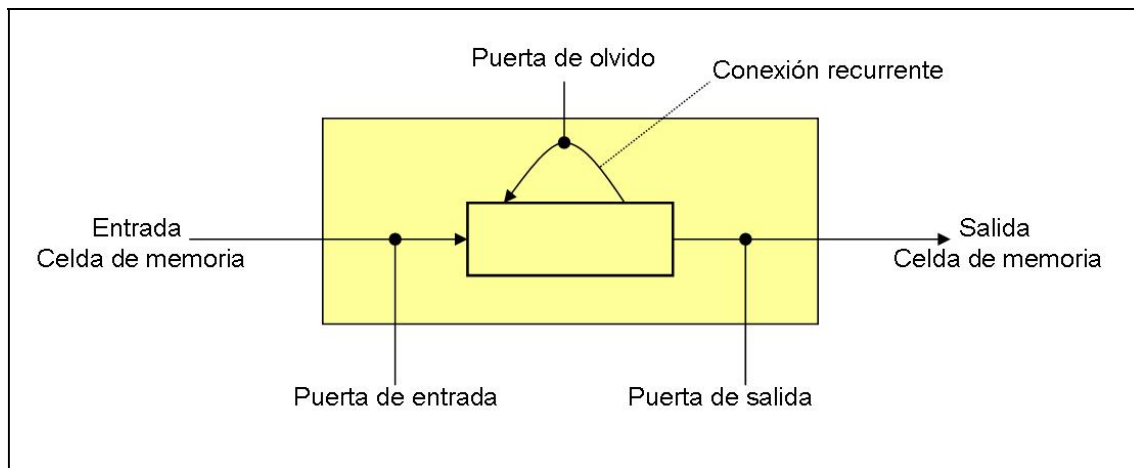


Figura 2.8 - Esquema de una neurona LSTM

Esta estructura, denominada *celda de memoria*, está compuesta de cuatro elementos: la *puerta de entrada*, una *conexión recurrente* a sí misma, una *puerta de olvido* y una *puerta de salida*. La *conexión recurrente* tiene un peso de 1 y asegura que, salvo una interferencia externa, el estado de la neurona puede permanecer constante entre las distintas iteraciones del entrenamiento. Las puertas sirven para modular las interacciones entre la propia celda y su entorno. La *puerta de entrada* permite que la señal recibida pueda alterar el estado de la celda o bloquearlo mientras que la *puerta de salida* permite que el estado se propague —o no— a otras neuronas. Por último, la *puerta de olvido* puede modular la conexión recurrente, determinando si la celda debe recordar u olvidar su estado previo.

El uso de redes LSTM se ha aplicado a diferentes problemas con un gran porcentaje de éxito: predicción de series temporales, reconocimiento del lenguaje, reconocimiento de imágenes y caracteres, etc.

3. Análisis

3.1 Introducción

En este apartado se incluye la colección de requisitos del sistema definidos al arranque del proyecto, con las directrices principales que debía cumplir la solución final implementada para la predicción de licitaciones de nuevos contratos.

A partir de esos requisitos y del conocimiento de las capacidades, alcance y restricciones de las herramientas a emplear, junto con el conocimiento adquirido durante la fase de análisis y estudio de los ficheros de entrada y los modelos más utilizados para tratamiento de series temporales, se ha elaborado el conjunto de requisitos software a satisfacer en las fases de diseño e implementación y que faciliten su verificación y completitud a lo largo del todo el proceso de desarrollo.

Se incluye también una matriz de trazabilidad entre los requisitos funcionales originales y los requisitos software.

3.2 Definición de requisitos de usuario

En este caso, ante la ausencia de cliente real, los requisitos de usuario han sido definidos en las sesiones de preparación del proyecto entre el tutor y el autor del proyecto. En ellos se recogen las líneas maestras que debe cumplir la aplicación para lograr el objetivo establecido, separando los requisitos en dos clases bien diferenciadas:

- **Requisitos de capacidad:** representan aquellos objetivos que se deben completar entendidos como las necesidades del usuario.
- **Requisitos de restricción:** formulan los límites que se deben considerar a la hora de abordar el proyecto, teniéndolos presentes a la hora de satisfacer los requisitos de capacidad.

Los atributos que se van a emplear para definir cada uno de los requisitos de usuario que se formulan a continuación son:

- **Identificador:** cada requisito está identificado de manera unívoca. El formato del identificador será de la forma UR-{C,R}-NN, siendo:
 - UR: prefijo que identifica al requisito como un requisito de usuario.
 - C, R: interfijo que indica si se trata de un requisito de capacidad –C- o de un requisito de restricción, R.
 - NN: secuencia numérica incremental que toma valores en el rango 01-99. Este secuencial será independiente para cada una de las categorías (capacidad, restricción).
- **Fuente:** informa del origen del requisito, atendiendo a quién lo ha formulado. En este caso va a tomar dos valores posibles: Tutor, Autor PFC.

- **Necesidad:** atendiendo a la obligatoriedad de cumplimentación de los requisitos se han definido tres niveles:
 - *Fundamental:* es obligatorio que el sistema final cumpla totalmente con lo especificado en el requisito.
 - *Recomendable:* sería deseable que el sistema final cumpliera total o parcialmente lo indicado en el requisito.
 - *Prescindible:* no supone ningún impacto en el usuario el hecho de no cumplir lo indicado en el requisito.
- **Prioridad:** se utilizan tres niveles de prioridad para poder priorizar y planificar la consecución de los requisitos. Por orden descendente de importancia los tres niveles de prioridad definidos son: *Alta, Media, Baja*.
- **Descripción:** cuerpo del requisito donde se expone el contenido, redactando el objetivo o restricción.

□ Requisitos de Capacidad

Identificador	UR-C-01		
Fuente	Tutor		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
Los ficheros a emplear como fuente de datos de contratación deben de ser los ficheros de texto históricos disponibles en la web de UsaSpending.gov.			

Tabla 3.1 - Requisito Usuario UR-C-01

Identificador	UR-C-02		
Fuente	Tutor		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema a desarrollar debe ser capaz de procesar de manera masiva los ficheros de datos de contratación y transformarlos en un resultado interpretable por modelos matemáticos o software que soporten la predicción de futuras licitaciones de contratos.			

Tabla 3.2 - Requisito Usuario UR-C-02

Identificador	UR-C-03		
Fuente	Autor PFC		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
Los ficheros de contratación disponen de hasta 250 indicadores sobre cada uno de los contratos. Se debe realizar un estudio de todos ellos para determinar cuáles son los más representativos de un contrato teniendo en cuenta el objetivo de predicción de una nueva licitación.			

Tabla 3.3 - Requisito Usuario UR-C-03

Identificador	UR-C-04		
Fuente	Tutor		
Necesidad	Fundamental	Prioridad	Media
Descripción			
Las transformaciones aplicadas sobre los datos origen deben garantizar la estabilidad y coherencia de los datos, sin alterar el orden o periodicidad de las series temporales que se puedan determinar en el transcurso del proyecto.			

Tabla 3.4 - Requisito Usuario UR-C-04

Identificador	UR-C-05		
Fuente	Tutor		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
Se debe experimentar con varios modelos predictivos y realizar una comparativa de los resultados obtenidos.			

Tabla 3.5 - Requisito Usuario UR-C-05

Identificador	UR-C-06		
Fuente	Tutor		
Necesidad	Fundamental	Prioridad	Baja
Descripción			
Se debe documentar todo el resultado del estudio previo y del proceso de desarrollo en una memoria de trabajo.			

Tabla 3.6 - Requisito Usuario UR-C-06

Identificador	UR-C-07		
Fuente	Autor PFC		
Necesidad	Recomendable	Prioridad	Baja
Descripción			
No se define ninguna interfaz gráfica determinada para la aplicación. La ejecución de los procesos se realizará en batch.			

Tabla 3.7 - Requisito Usuario UR-C-07

□ Requisitos de Restricción

Identificador	UR-R-01		
Fuente	Tutor		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
Se debe emplear Apache Flink como librería para el tratamiento masivo de los ficheros de entrada.			

Tabla 3.8 - Requisito Usuario UR-R-01

Identificador	UR-R-02		
Fuente	Tutor		
Necesidad	Prescindible	Prioridad	Baja
Descripción			
Se recomienda que una de las soluciones empleadas para la predicción de observaciones futuras esté basada en la herramienta TensorFlow.			

Tabla 3.9 - Requisito Usuario UR-R-02

Identificador	UR-R-03		
Fuente	Autor PFC		
Necesidad	Fundamental	Prioridad	Media
Descripción			
El entorno de desarrollo y ejecución de la solución a desarrollar será el propio ordenador del autor del proyecto.			

Tabla 3.10 - Requisito Usuario UR-R-03

3.3 Establecimiento de los requisitos software

A partir de los requisitos de usuario definidos en el apartado anterior y habiendo realizado un estudio previo de los ficheros de datos y de las herramientas a utilizar a lo largo del proyecto se han elaborado los requisitos software.

Los atributos para representarlos son:

- **Identificador:** cada requisito está identificado de manera unívoca. El formato del identificador será de la forma SR-{Clasificación del requisito software}-NN, siendo:

- SR: prefijo que identifica al requisito como un requisito de software.
- Clasificación del requisito software: interfijo que identifica a que conjunto, de los detallados posteriormente, pertenece este requisito.

- NN: secuencia numérica incremental que toma valores en el rango 01-99. Este secuencial será independiente para cada una de las categorías.
- **Fuente:** informa del origen del requisito, atendiendo con cuál o cuáles requisitos de usuario y/o software está relacionado.
- Los campos de **Necesidad** y **Prioridad** se mantienen en significado y forma igual que para los requisitos de usuario.

Las categorías contempladas en este proyecto para la clasificación de los requisitos software han sido:

- **Requisitos Funcionales** (interfijo F): aquellos relativos a qué tiene que hacer el software. Se derivan de los requisitos de capacidad del usuario (puesto que sólo se define un caso de uso: procesar ficheros y modelar para predicción).
- **Requisitos de Operación** (interfijo O): especifican cómo va a alcanzar la aplicación la funcionalidad requerida.
- **Requisitos de Recursos** (interfijo R): especifican las necesidades hardware y software de terceros que debe satisfacer la solución.
- **Requisitos de Documentación** (interfijo D): se indica qué documentación asociada al proyecto debe elaborarse.

No se han definido requisitos relativos a Calidad, Rendimiento, Interfaz con sistemas terceros, Seguridad o Disponibilidad del sistema, por eso no se incluyen estas categorías habituales entre las empleadas en la clasificación de requisitos software de este proyecto.

□ Requisitos Funcionales

Identificador	SR-F-01		
Fuente	UR-C-01		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema empleará como datos de partida aquellos ficheros de texto plano con extensión .csv descargados desde la web de UsaSpending.gov para el periodo 2000-2015. Contienen una primera línea de cabecera. Cada línea representa información de un contrato y se dispone de 250 características, separadas entre sí por una coma ','. Cada característica aparece entrecomillada doblemente.			

Tabla 3.11 - Requisito Software SR-F-01

Identificador	SR-F-02		
Fuente	UR-C-02		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema realizará una primera transformación a los ficheros de datos originales para reducirlos a aquellas características que pueden ser mejores indicadores para los modelos de predicción y para modificar su formato, deshaciendo el doble entrecomillado.			

Tabla 3.12 - Requisito Software SR-F-02

Identificador	SR-F-03		
Fuente	UR-C-03		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
Se llevará a cabo un análisis pormenorizado de cada una de las características preseleccionadas para determinar si se usan como identificadores de un contrato.			

Tabla 3.13 - Requisito Software SR-F-03

Identificador	SR-F-04		
Fuente	UR-C-02, UR-C-04		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema sólo tratará de determinar el comportamiento futuro de nuevas licitaciones de contratos, prescindiendo de aquellas líneas de los ficheros de contratos que estén relacionadas con modificaciones de contratos ya adjudicados.			

Tabla 3.14- Requisito Software SR-F-04

Identificador	SR-F-05		
Fuente	UR-C-02, UR-C-03		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema agrupará el número de contratos licitados en función de las características escogidas y del mes y año en que se ha adjudicado, construyendo la serie temporal de adjudicaciones con periodicidad mensual para el rango de años escogido.			

Tabla 3.15- Requisito Software SR-F-05

Identificador	SR-F-06		
Fuente	UR-C-02, UR-C-04		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema empleará el 0 para aquellos elementos de la serie temporal determinados por un par mes-año en el que se carezca de contratos adjudicados para unas determinadas características.			

Tabla 3.16- Requisito Software SR-F-06

Identificador	SR-F-07		
Fuente	UR-C-05		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema tratará de modelar el problema empleando regresión múltiple lineal, empleando distintas características de los contratos como predictores y evaluando la idoneidad del modelo obtenido.			

Tabla 3.17- Requisito Software SR-F-07

Identificador	SR-F-08		
Fuente	UR-C-05		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema empleará un modelo ARIMA para las series temporales de licitaciones de contratos y se evaluarán los resultados obtenidos.			

Tabla 3.18- Requisito Software SR-F-08

Identificador	SR-F-09		
Fuente	UR-C-05		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema implementará una red de neuronas LSTM para la predicción de futuras observaciones. Se deben probar distintas configuraciones de la red y elegir una de ellas para evaluar los resultados obtenidos.			

Tabla 3.19- Requisito Software SR-F-09

Identificador	SR-F-10		
Fuente	UR-C-05		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
Se compararán los resultados obtenidos por los distintos modelos de predicción aplicados por el sistema, empleando como método la raíz cuadrada del error cuadrático medio obtenido entre los datos reales y los obtenidos por cada modelo. Se elaborarán las conclusiones oportunas.			

Tabla 3.20- Requisito Software SR-F-10

❑ Requisitos de Operación

Identificador	SR-O-01		
Fuente	UR-R-01		
Necesidad	Fundamental	Prioridad	Alta
Descripción			
El sistema empleará Java utilizando la plataforma de Apache Flink para el procesamiento de los ficheros de texto y su transformación a ficheros con series temporales.			

Tabla 3.21- Requisito Software SR-O-01

Identificador	SR-O-02		
Fuente	UR-C-03		
Necesidad	Recomendable	Prioridad	Media
Descripción			
El sistema empleará el lenguaje R para las siguientes tareas: Para el análisis de características de los contratos, en cuanto al estudio del dominio donde toman valores y frecuencia de repetición de los mismos. Para la aplicación del modelo de regresión múltiple lineal. Para la aplicación del modelo ARIMA mediante la función <i>auto.arima</i> .			

Tabla 3.22- Requisito Software SR-O-02

Identificador	SR-O-03		
Fuente	UR-R-02		
Necesidad	Prescindible	Prioridad	Baja
Descripción			
Se emplea Docker como solución para la instalación y ejecución de Keras y TensorFlow, utilizadas como librería y motor de ejecución de una red de neuronas LSTM.			

Tabla 3.23- Requisito Software SR-O-03

Identificador	SR-O-04		
Fuente	UR-C-07		
Necesidad	Recomendable	Prioridad	Baja
Descripción			
Para facilitar la ejecución de la aplicación de transformación de los datos de entrada, las clases Java tienen asociados unos ficheros de configuración donde se especifican los parámetros necesarios para su ejecución, siendo mayoritariamente éstos las rutas de entrada y salida de los ficheros transformados.			

Tabla 3.24- Requisito Software SR-O-04

☐ Requisitos de Recursos

Identificador	SR-R-01		
Fuente	UR-R-03		
Necesidad	Fundamental	Prioridad	Media
Descripción			
Las características del equipo a utilizar como entorno de desarrollo y ejecución de la aplicación son: <i>Equipo: Toshiba Satellite Pro C70-B-34T</i> Procesador: <i>Intel Core i5-5200U 2.20GH</i> . 4 procesadores lógicos. Memoria RAM: 16 GB Sistema Operativo: Windows 10 Pro V.1511			

Tabla 3.25- Requisito Software SR-R-01

☐ Requisitos de Documentación

Identificador	SR-D-01		
Fuente	UR-C-06		
Necesidad	Fundamental	Prioridad	Baja
Descripción			
Se elaborará una memoria de trabajo donde se incluya un estado del arte de las tecnologías y áreas de conocimiento que aplican al objeto del proyecto, un resumen de los productos obtenidos a lo largo del ciclo de vida del desarrollo del proyecto y una descripción de la implementación y los resultados obtenidos durante la experimentación, elaborando una conclusiones y una lista de sugerencias de trabajo futuras que puedan aportar mejoras o alternativas al proyecto elaborado.			

Tabla 3.26- Requisito Software SR-D-01

Identificador	SR-D-02		
Fuente	UR-C-06		
Necesidad	Fundamental	Prioridad	Baja
Descripción			
Se elaborará un manual de usuario donde se expliquen los ficheros de configuración y los parámetros contenidos en éstos necesarios para la correcta ejecución de la aplicación desarrollada.			

Tabla 3.27- Requisito Software SR-D-02

3.4 Matriz trazabilidad requisitos usuario/SW

	UR-C-01	UR-C-02	UR-C-03	UR-C-04	UR-C-05	UR-C-06	UR-C-07	UR-R-01	UR-R-02	UR-R-03
SR-F-01	X									
SR-F-02		X								
SR-F-03			X							
SR-F-04		X	X	X						
SR-F-05		X	X	X						
SR-F-06		X	X	X						
SR-F-07					X	X				
SR-F-08					X	X				
SR-F-09					X	X				
SR-F-10					X	X				
SR-O-01								X		
SR-O-02			X							
SR-O-03									X	
SR-O-04							X			
SR-R-01										X
SR-D-01						X	X			
SR-D-02						X	X			

Figura 3.1 - Matriz trazabilidad Req. Usuario/SW

4. Diseño

4.1 Arquitectura de la solución. Diagrama de Clases.

Para dar solución al problema este se ha dividido en dos partes diferenciadas:

- Primero, el tratamiento de los ficheros de datos de los contratos hasta conseguir transformarlos en un formato válido para poderles aplicar métodos matemáticos que permitan la generación de un modelo que facilite la predicción. Las 6 fases que se engloban dentro se representan en la secuencia de la Figura 4.1. Para aquellas fases de esta figura donde aparece el logotipo de Apache Flink y Java se han implementado las clases necesarias para implementar la funcionalidad requerida a ejecutar en ella: el diagrama de clases de todas las empleadas en esta parte de la implementación se representa en la Figura 4.3.

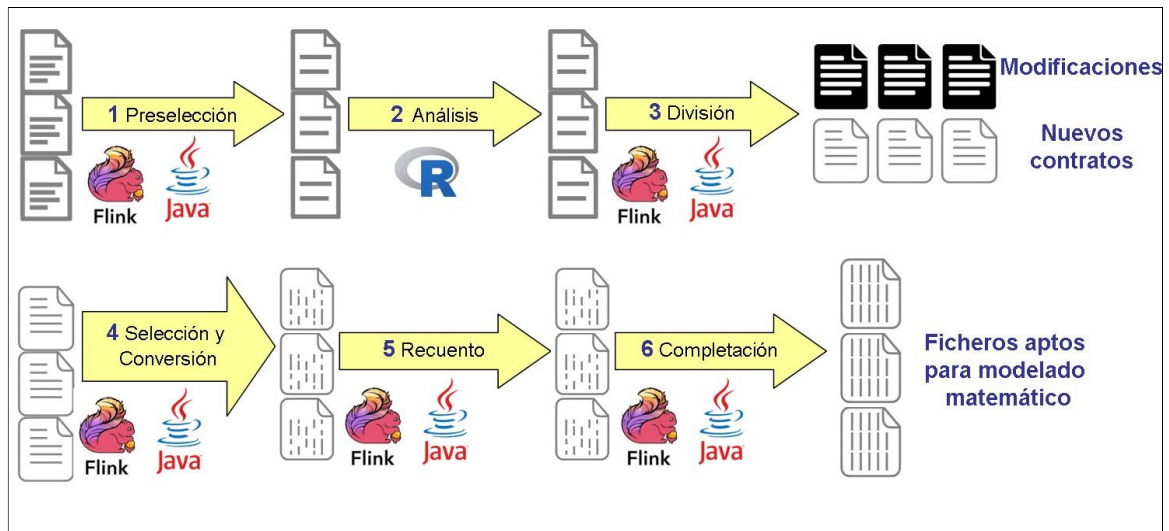


Figura 4.1 -Secuencia para el tratamiento de los ficheros

- Segundo, una vez elaborados los ficheros, se les aplican distintos métodos de aproximación para el modelado de las series temporales y se comparan los resultados obtenidos. En la Figura 4.2 se representa la secuencia de tareas que se detallan en el texto más adelante.

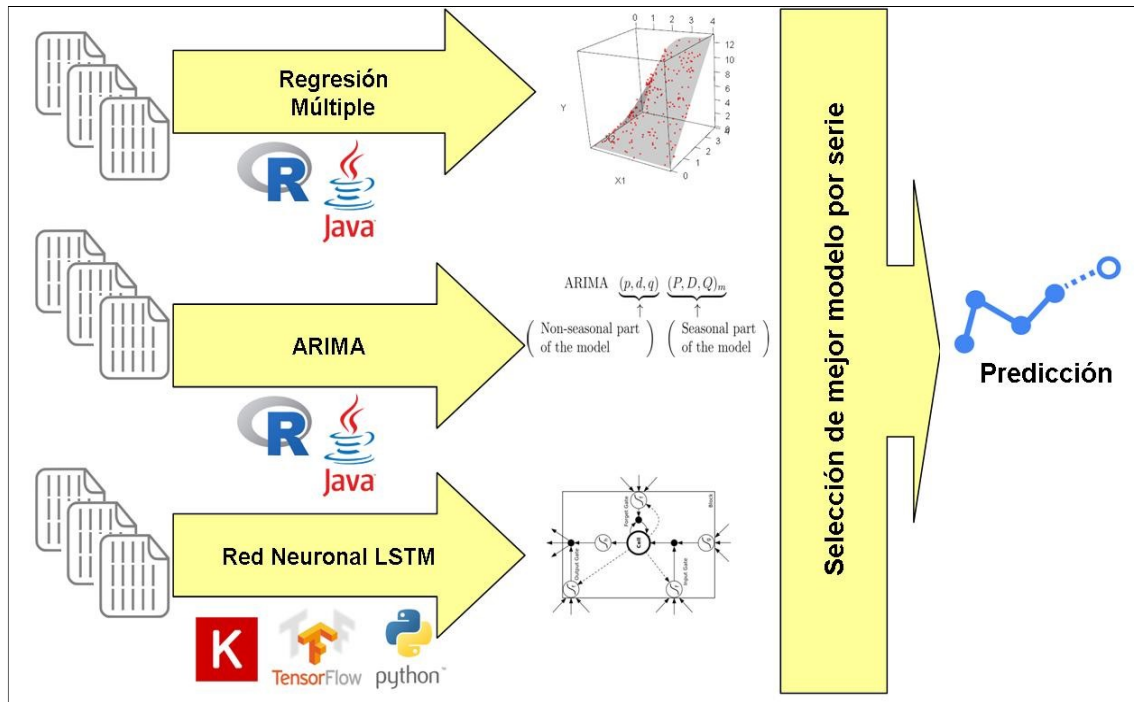


Figura 4.2 - Aplicación de modelos y selección del más adecuado

La aplicación de la regresión múltiple lineal y de los modelos ARIMA se ha realizado utilizando R como lenguaje. Para la selección del mejor modelo, detallado en el apartado de Experimentación, se ha desarrollado una nueva clase Java que utiliza R desde ella empleando JRI.

La red de neuronas artificial se ha implementado utilizando Keras sobre TensorFlow, empleando como lenguaje Python.

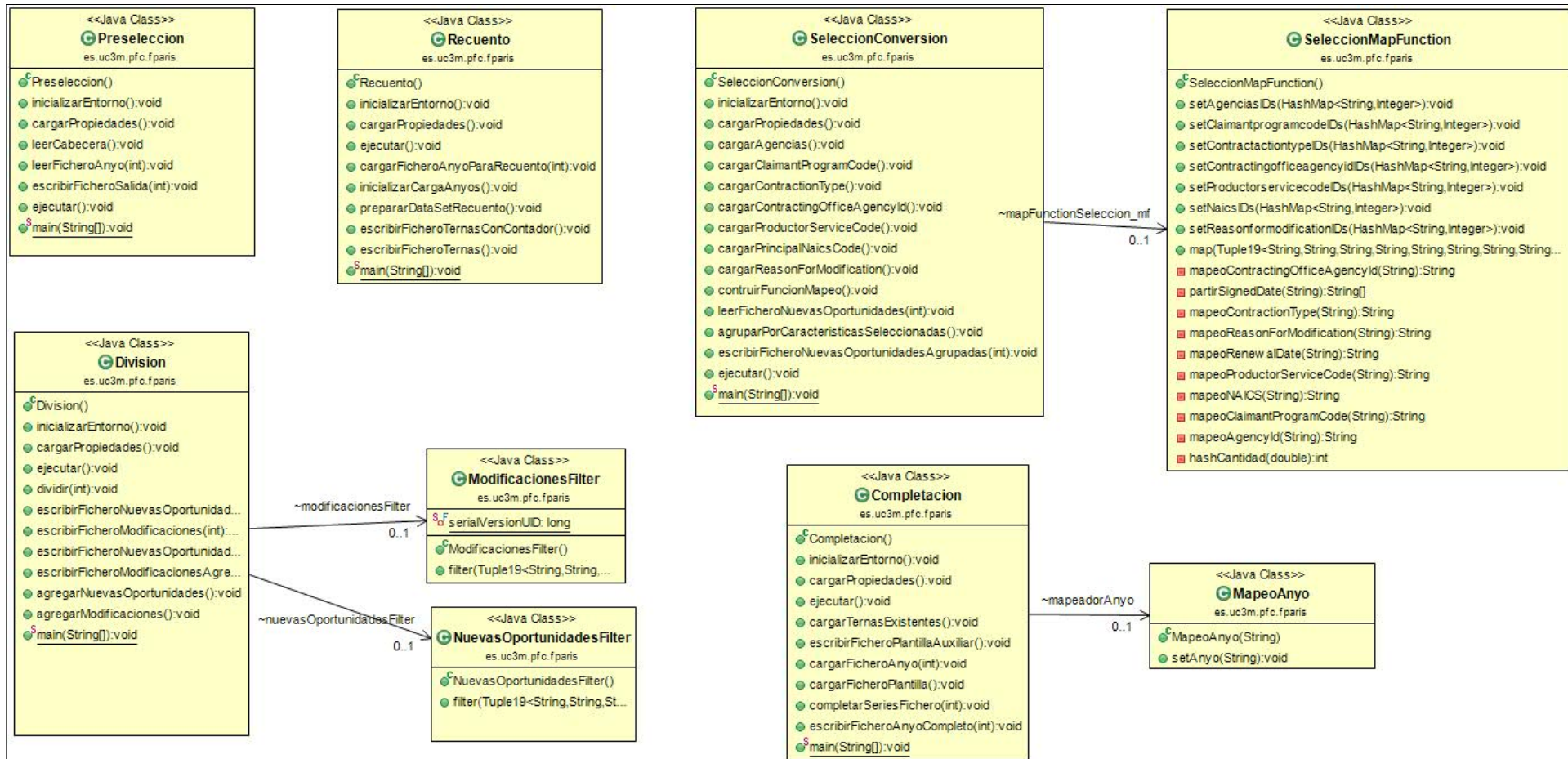


Figura 4.3 - Diagrama de Clases empleadas en la fase de Implementación

4.2 Entorno tecnológico

4.2.1 Apache Flink

En paralelo al desarrollo de Apache Hadoop y Spark, en 2010 el proyecto de investigación "*Stratosphere: Information Management on the Cloud*" [9] en colaboración con distintas universidades alemanas conciben lo que en 2014 se acaba convirtiendo en Apache Flink [25].



Figura 4.4 - Logotipo de Apache Flink

Para entender mejor las distintas funcionalidades que ofrece Apache Flink lo mejor es empezar presentando el conjunto de componentes del que dispone. En la Figura 4.5 se muestra el *stack* de Flink.

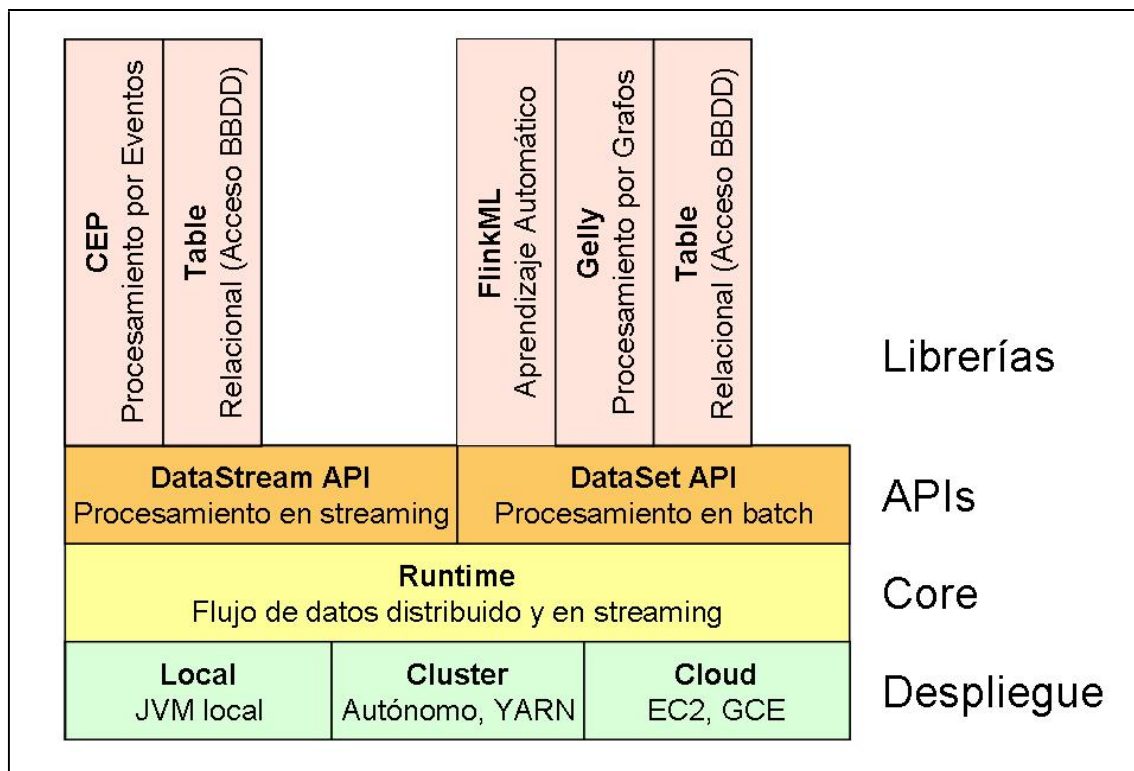


Figura 4.5 - Arquitectura y componentes de Apache Flink

Flink se puede desplegar de manera autónoma, sin necesidad de ningún otro componente del entorno de Hadoop aunque es habitual que para su despliegue en clúster se apoye en YARN -de Hadoop- para la gestión de recursos y que se utilice Flink

combinado con HDFS, el gestor de ficheros de Hadoop. El despliegue puede hacerse a nivel local, en clúster o en la *nube*.

Flink permite el procesamiento tanto en *batch* como en *streaming*. Puesto que el alcance del proyecto queda dentro del procesamiento *batch* se explica este en más detalle.

Se puede programar tanto en Java como en Scala utilizando las librerías de Flink. El objeto básico para el procesamiento en *batch* es un *DataSet* que representa un conjunto de datos (*DataStream* para el procesamiento en *streaming*). Sobre cada *dataset* se aplican una serie de transformaciones que originan nuevos *datasets* de salida. El origen de datos de cada uno de esos *dataset* se le denomina *Source* y a la salida o uso que se le da se denomina *Sink*.

Todos los programas por defecto en Flink son paralelos y distribuidos. Los *datasets* se dividen en particiones a las que se aplican las transformaciones (subtareas) en tantos hilos (*threads*) o servidores como se hayan configurado. Al final se aplican las transformaciones necesarias para reunir esas particiones del *dataset* y originar la salida correspondiente. Los objetos de los *dataset* deben ser *serializables* para que se pueda aplicar el paralelismo.

En cada ejecución de Flink se despliegan dos tipos de procesos:

- ***JobManager***: es un proceso que actúa como *master* (líder), programando las subtareas y coordinando la comunicación entre ellas. Se encarga también de establecer puntos de control de la ejecución y de la recuperación ante fallos. Como mínimo siempre hay un *JobManager*, incluso en ejecuciones en local. En ejecuciones distribuidas es típico tener un conjunto de *JobManagers*, entre los cuáles uno es el líder.
- ***TaskManager***: son los procesos que se encargan de ejecutar las subtareas de transformación sobre las particiones de los *dataset* y del intercambio de información con otros *taskmanagers* y con el *JobManager*. Cada *TaskManager* emplea una máquina virtual de Java (*JVM*, del inglés, *Java Virtual Machine*) independiente si se ejecuta en modo distribuido; en local comparten una misma JVM el *JobManager* y el *TaskManager*.

La estructura básica de un programa en Flink consta de las siguientes partes:

- 1 Obtener un entorno de ejecución. Decide dinámicamente si se está ejecutando en local o en clúster.
- 2 Carga de los datos.
- 3 Indicar las transformaciones de los datos.
- 4 Indicar la salida de los datos transformados.
- 5 Lanzar la ejecución del programa.

En la Figura 4.6 puede verse un ejemplo real de un programa Flink correspondiente a algunas líneas de una clase Java de la implementación de este proyecto y un esquema de la serie de transformaciones aplicadas, identificando cada uno de los elementos.

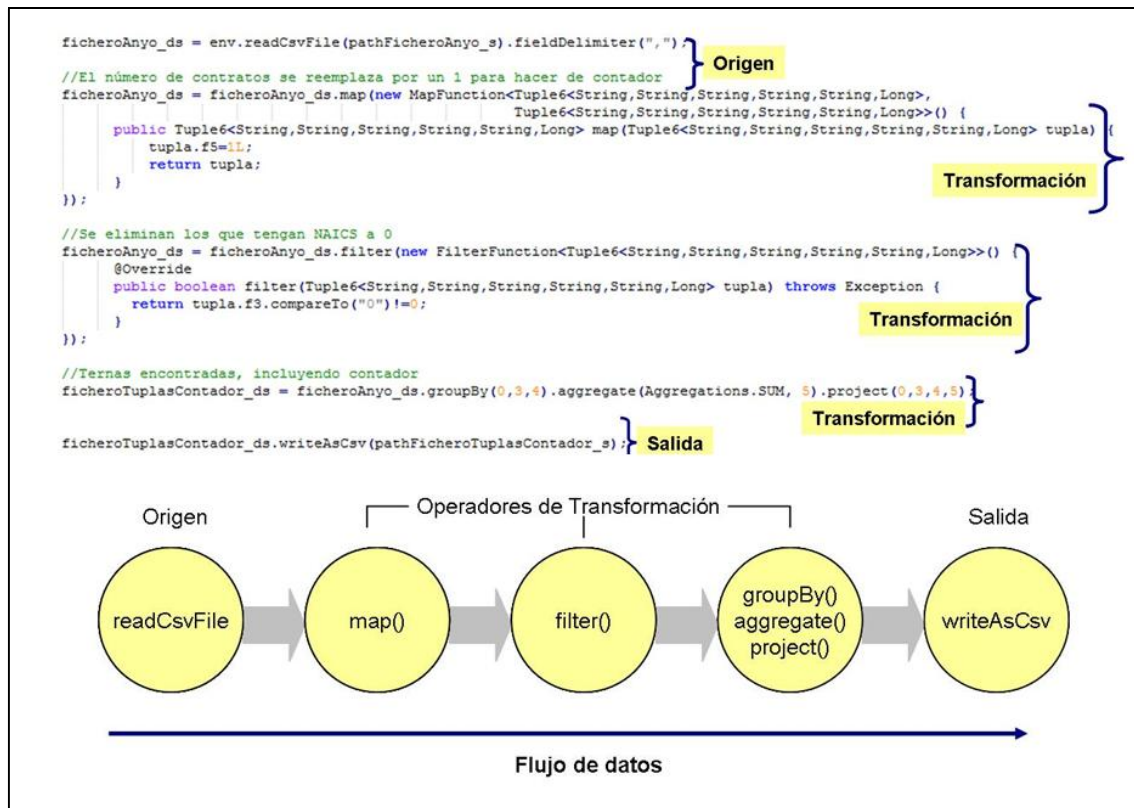


Figura 4.6 - Ejemplo programa en Apache Flink y cadena de transformaciones

La ejecución del programa en Flink no se realiza directamente una vez se ejecuta el método *main*. Cada operación de carga de datos, transformación y salida se añade al plan de ejecución y éste se reordena automáticamente con el fin de optimizar la ejecución del proceso. Una vez diseñado el plan de ejecución se realiza la ejecución propiamente dicha. La capa de ejecución recibe un programa en la forma de *JobGraph*. Un *JobGraph* es un flujo genérico de datos a ejecutarse en paralelo con tareas que producen y consume datos. Las dos APIs (*DataStream* y *DataSet API*) generan estos *JobGraphs* mediante procesos de compilación independientes. El API para *DataSet* usa un optimizador para elaborar el mejor plan de ejecución para el programa, mientras que el API para *DataStream* usa un constructor para *streams* (flujos).

Flink ofrece una interfaz web (*Apache Flink Web Client*) que actúa como interfaz gráfica cliente para permitir la gestión de los procesos a ejecutar, incorporando ficheros *JARs* para su ejecución. Una de las mayores ventajas que ofrece ejecutar los programas desde esta herramienta es que se puede visualizar el plan de ejecución optimizado para ese programa antes de ejecutarse, como se muestra en la Figura 4.7.

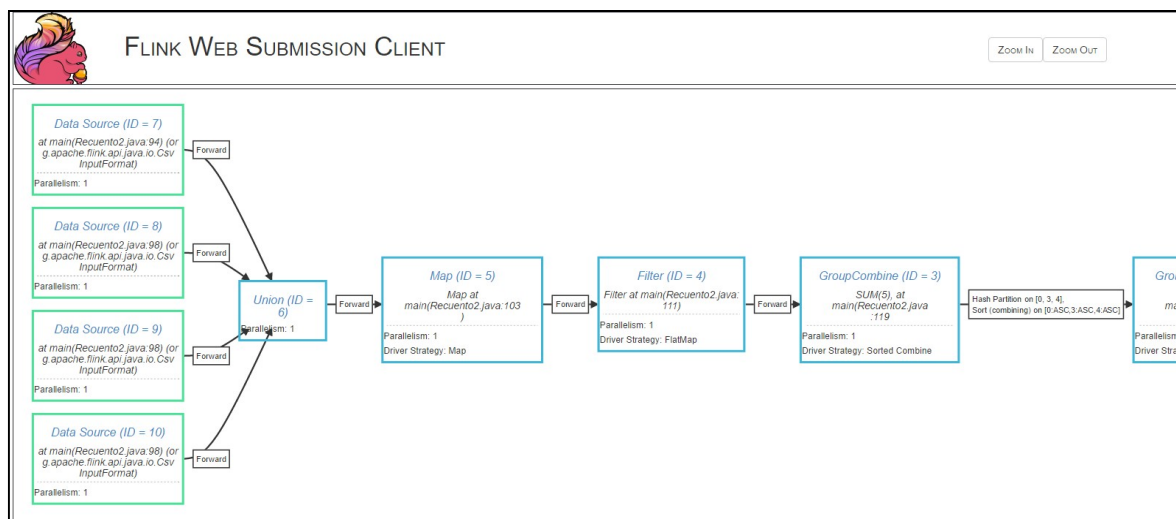


Figura 4.7 - Vista de la interfaz web cliente con un mapa optimizado de ejecución

Además de esta interfaz cliente que permite la carga y ejecución de aplicaciones, existe otra consola que permite comprobar la ejecución de todos los trabajos: *Apache Flink Dashboard*. Desde este cuadro de mandos se puede comprobar el grado de avance de una tarea, comprobar su resultado o ver el tiempo que se ha invertido en ejecutar cada una de las transformaciones identificadas en el plan de ejecución. En la Figura 4.8 se muestra el resultado de las tareas ejecutadas –incluyendo nombre, tiempo empleado y estado de la ejecución, entre otros- que aparecían en el plan de ejecución de la Figura 4.7 que se mostraba anteriormente.

Flink Java Job at Wed May 18 13:14:45 CEST 2016 32a2171bdb26b6d85da4faa54ebb47dc 0 0 0 0 0 0 0 0 0 0 2016-05-18, 13:14:46 - 2016-05-18, 13:15:01									
Plan Timeline Exceptions Properties Configuration									
Overview Accumulators									
Start Time	End Time	Duration	Name	Bytes received	Records received	Bytes sent	Records sent	Tasks	Status
2016-05-18, 13:14:46	2016-05-18, 13:14:47	627ms	DataSource (at main(Recuento2.java:98) (org.apache.flink.api.java.io.CsvInputFormat))	0	0	5,780,914	221,012	0 0 0 1	FINISHED
2016-05-18, 13:14:46	2016-05-18, 13:14:47	643ms	DataSource (at main(Recuento2.java:98) (org.apache.flink.api.java.io.CsvInputFormat))	0	0	5,679,938	217,046	0 0 0 1	FINISHED
2016-05-18, 13:14:46	2016-05-18, 13:14:47	643ms	DataSource (at main(Recuento2.java:94) (org.apache.flink.api.java.io.CsvInputFormat))	0	0	6,434,334	245,771	0 0 0 1	FINISHED
2016-05-18, 13:14:46	2016-05-18, 13:14:47	627ms	DataSource (at main(Recuento2.java:98) (org.apache.flink.api.java.io.CsvInputFormat))	0	0	5,843,072	223,208	0 0 0 1	FINISHED
2016-05-18, 13:14:46	2016-05-18, 13:14:59	12s	CHAIN Map (Map at main(Recuento2.java:103)) -> Filter (Filter at main(Recuento2.java:111)) -> Combine(SUM(5), at main(Recuento2.java:119))	23,738,258	907,037	3,629,372	138,935	0 0 0 1	FINISHED
2016-05-18, 13:14:59	2016-05-18, 13:15:01	2s	CHAIN Reduce (SUM(5), at main(Recuento2.java:119)) -> Map (Projection [0, 3, 4, 5])	3,629,372	138,935	5,270,398	277,870	0 0 0 1	FINISHED
2016-05-18, 13:15:01	2016-05-18, 13:15:01	601ms	DataSink (CsvOutputFormat (path: C:/PFC/workspace/PruebaMemoria/tuplas_encontradas_Contador.txt, delimiter: .))	2,635,199	138,935	0	0	0 0 0 1	FINISHED
2016-05-18, 13:15:01	2016-05-18, 13:15:01	585ms	Map (Projection [0, 1, 2])	2,635,199	138,935	1,523,719	138,935	0 0 0 1	FINISHED
2016-05-18, 13:15:01	2016-05-18, 13:15:01	586ms	DataSink (CsvOutputFormat (path: C:/PFC/workspace/PruebaMemoria/tuplas_encontradas.txt, delimiter: .))	1,523,719	138,935	0	0	0 0 0 1	FINISHED

Figura 4.8 - Vista del cuadro de mando con el resultado de ejecución

Para arrancar esta consola es necesario hacerlo desde los scripts de Apache Flink, ejecutando *start-local.sh*. Por defecto está accesible en la siguiente URL: <http://localhost:8081>.

En la Figura 4.9 se incluye una imagen de la vista general del cuadro de mandos cuando se accede a él. Se puede ver a la izquierda el menú de opciones para la ejecución y el control de los trabajos.

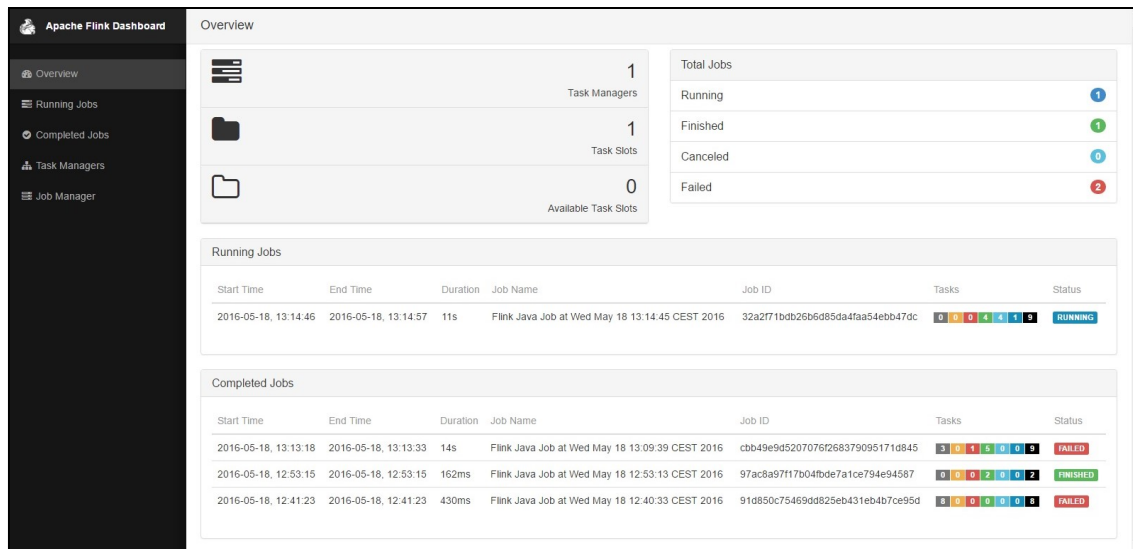


Figura 4.9 - Vista principal del cuadro de mandos

□ Fuentes, salidas y transformaciones más habituales

Flink no se limita sólo a incluir transformaciones *MapReduce* sobre los ficheros de entrada. Permite la importación de datos desde diferentes orígenes, su volcado en diferentes formas y soportes y la aplicación de una rica colección de transformaciones. Para cada uno de los elementos estas son las posibilidades más habituales:

Fuentes (Sources)

- **Desde fichero:** permite la lectura de datos de ficheros con distintos formatos predefinidos (txt, csv, etc.) o indicar el formato personalizado del fichero antes de cargarlo.

Para un procesamiento en *streaming* se puede incluso utilizar *readFileStream* que va creando los bloques a partir de los cambios que se van incorporando al fichero.

- **Desde sockets:** *socketTextStream* permite la lectura de elementos separados por un delimitador utilizando *sockets*.

- **Desde un origen de datos personalizado:** permite incorporar cualquier origen que suministre una colección de datos. Por ejemplo, con *addSource* se puede incluir un consumidor que suministra datos a partir de Apache Kafka [39] (utilizado para colas de mensajería).

- **Desde una colección de objetos:** con el método *fromCollection* se puede crear un *dataset* a partir de cualquier objeto Java que implemente la interfaz *java.util.Collection*. Con *fromElements*, pasándole como parámetro una secuencia de elementos crea un *dataset* a partir de ellos. Una limitación para estas dos formas

de construir un origen de datos es que todos los elementos de la colección o de la lista recibida como parámetro deben ser del mismo tipo.

Salidas (Sinks)

Prácticamente cada tipo de origen de datos tiene su correspondencia en su tipo para constituir una salida que consuma *datasets*:

- **Escritura a fichero:** con formatos predefinidos (*writeAsCsv*) o personalizados (*writeAsFormattedText*) donde se indica los elementos a incorporar a la salida y sus delimitadores.
- **Volcado a salida estándar:** con *print* o *printToErr* se vuelca el contenido a la salida estándar.
- **Volcado a sockets.**
- **Salidas personalizadas:** con *addSink* se puede incluir el conector que se necesite entre los implementados. Por ejemplo, de igual manera que se consumía con Apache Kafka de la cola de mensajes se puede también incluir mensajes en ella.

Es importante tener presente que Flink considera como *sink* cualquier uso que se haga sobre el *dataset* y que genere un resultado de salida que se pueda considerar producto de la entrada. Por ejemplo, si sobre el *dataset* se aplican las escrituras a fichero, socket, o salida personalizada, o si se envía el contenido a la salida estándar, se va a considerar que se ha empleado el *dataset*. Pero también se va a considerar utilizado si a partir de él se genera una colección de objetos (aplicando el método *collect* sobre el *dataset*) o si se cuenta el número de elementos que lo componen (utilizando el método *count*). Cualquiera de estos métodos que implica un uso de un *dataset* que termina en una salida ya es considerado por Flink como un bloque de ejecución. Cuando desde el entorno de ejecución se hace uso de la función *execute* y no se han definido nuevos bloques de programa tras la última generación de una salida (*sink*) se genera una excepción de este tipo:

```
java.lang.RuntimeException: No new data sinks have been defined since the last execution. The last execution refers to the latest call to 'execute()', 'count()', 'collect()', or 'print()'.
```

Transformaciones

Las transformaciones más habituales que se pueden aplicar a los elementos de un *dataset* y que han sido empleadas en la fase de implementación son:

- **Map:** aplica un mapeo uno-a-uno a cada uno de los elementos de un *dataset*: a cada elemento de entrada se le asigna un elemento de salida.

- **FlatMap:** similar a la transformación *Map* pero permitiendo que a cada elemento de entrada se le mapee a uno, varios o incluso ningún elemento de salida.
- **Filter:** aplica un filtro mediante una función que implemente la interfaz *FilterFunction* a los datos del *dataset* de entrada dejando sólo aquellos que cumplen con las condiciones que exige la función. Esta función debe devolver *true* si se cumplen las condiciones y *false* en caso contrario.
- **Project:** está disponible sólo para la interfaz Java de Flink y es aplicable sólo para *datasets* formados por tuplas que permiten el indizado numérico de los campos. Lo que facilita la transformación *Project* es la selección de los campos de una tupla que se quieren mantener en la salida de la transformación desechando aquellos que no se incluyen como parámetro.
- **Reduce:** se aplica una función de reducción a todos los elementos de un conjunto de datos combinando iterativamente pares de elementos en uno solo hasta que quede un único en la salida.
- **GroupReduce:** primero es necesario agrupar los elementos de un *dataset* mediante un criterio de agrupación con la función *groupBy*. A continuación se aplica una función *GroupReduceFunction* a todo el grupo ordenado de elementos a la vez. A diferencia de *reduce* no tiene porqué devolver un único elemento de salida.
- **Agregaciones:** una vez que se han agrupado todos los elementos de un conjunto de datos con *groupBy* se permite la aplicación de una agregación entre Suma, Máximo o Mínimo. En un futuro está previsto que se extienda la funcionalidad de las agregaciones, puesto que ahora ofrecen un número muy limitado de ellas (por ejemplo, no es posible hacer un *count* del número de elementos si no es sumando un campo, ni tampoco utilizar para agregar campos que no sean tipos básicos: si se pretende, por ejemplo, agregar por un campo que sea un *BigInteger* se genera una excepción que indica que ese tipo de objeto no está soportado).
- **Ordenaciones:** sobre un *dataset*, agrupado o no, se pueden solicitar ordenaciones por diferentes criterios bien mediante *sortPartition* si no está agrupado o bien mediante *sortGroup*.
- **Union:** similar a *join* pero sin necesidad de indicar ningún criterio de unión: se unen todos los elementos de ambos *datasets* en uno sólo.

Existen algunas otras, como **Join** (permite la combinación de dos conjuntos de datos indicando un criterio de unión o **Cross**, que pueden consultarse en [41] aunque no se han empleado en ninguna de las clases Java de la implementación.

□ Aprendizaje automático con Apache Flink

Apache Flink dispone de **FlinkML**, una librería independiente al *core* y que implementa algunas de las utilidades más habituales para aprendizaje automático (el sufijo *ML* de FlinkML hace referencia a *Machine Learning*). A día de hoy está en versión beta (0.10.1) y en pleno desarrollo; se prevé ir completando con nuevas utilidades en el futuro y puede consultarse el plan de evolución de la misma en [12]. Sólo está disponible a día de hoy para programación con Scala.

Las funcionalidades implementadas hasta el momento son:

Para aprendizaje supervisado:

- Regresión lineal múltiple.
- SVM (*Support Vector Machines*) usando el algoritmo CoCoA (*Communication-efficient distributed dual Coordinate Ascent*), utilizado típicamente en problemas de clasificación y regresión.
- Optimización, que implementa métodos como el SGD (*Stochastic Gradient Descent*) o Funciones de Pérdida (*Loss Functions*).

Preprocesado de datos:

- Transformador de características polinómico.
- *Standard Scaler*: escalado de los datos de entrada atendiendo a la media y a la varianza.
- *MinMax Scaler*: para escalar los datos de entrada entre un rango [máximo, mínimo] definido por el usuario.

Recomendación:

- Algoritmo *Alternating Least Squares* (ALS).

Utilidades:

- Métricas para el cálculo de distancias según diferentes criterios (distancia euclídea, distancia Manhattan, etc.).

De momento hay algunas funcionalidades muy populares en otras soluciones de aprendizaje automático como la transformación de características continuas en valores discretos, aprendizaje no supervisado –*clustering*, por ejemplo–, reducción de las dimensiones de las características de entrada, árboles de decisión, etc., que no han sido

aún implementadas pero que forman parte del trabajo futuro a implementar dentro de la librería.

□ Elección de Flink. Comparativa Spark-Flink

Queda fuera del alcance de este proyecto establecer una comparativa en profundidad sobre las distintas soluciones *Big Data* presentadas para tratar grandes volúmenes de datos. Hadoop, Spark, Storm [40], Flink... cada uno presentan unas ventajas e inconvenientes y para el tratamiento de los ficheros que en este caso se requiere cualquiera de las soluciones podría haberse aplicado.

La elección de Apache Flink se ha basado en que es la solución más reciente y quizás la menos extendida hasta este momento, así que se consideraba interesante testarla dentro de este proyecto y superar la dificultad de disponer de menos documentación y una comunidad de usuarios más reducida.

No obstante, es importante señalar que uno de los puntos fuertes de Flink frente a Spark es la gestión de memoria dentro de la JVM. Spark puede serializar datos a disco pero siempre requiere que una parte de ellos se encuentre en memoria (concretamente en el *heap* de la JVM) para determinadas operaciones. Si el tamaño destinado a esta parte de la memoria no es suficiente, el programa en ejecución finaliza con una excepción por falta de memoria (*OutOfMemoryError*). Por el contrario, Flink está diseñado para no acumular grandes cantidades de objetos en el *heap* de la máquina virtual y sí hacerlo en una región de memoria distinta diseñada para tal fin. Todas las transformaciones han sido diseñadas para consumir el mínimo de memoria y poder volcar a disco cuando sea necesario.

Apache Spark ha anunciado un nuevo proyecto bautizado como *Tungsten* que promete eliminar esa sobrecarga en el *heap* de la JVM [10].

□ Otras funcionalidades de Apache Flink

En la Figura 4.5 se mostraban la arquitectura de Apache Flink. Dentro de las librerías disponibles, además de la de FlinkML para aprendizaje automático aparecen otras especializadas en:

- **FlinkCEP** [45]: es la librería dedicada al procesamiento por eventos en *streaming*. Se considera evento la ocurrencia de un patrón dentro del flujo de datos. Cada patrón se compone de varias etapas o estados. Con el fin de pasar de un estado al siguiente, el usuario puede especificar las condiciones de transición: contigüidad de eventos, filtros que deben cumplirse, etc.
- **Table API y SQL** [47]: este API, disponible para *batch* y *streaming*, ofrece una abstracción de la tabla relacional sobre los *datasets* que se va a procesar. Se pueden aplicar los operadores relacionales, tales como la selección, agregación y unión entre las tablas. Para consultar una tabla utilizando sintaxis SQL es necesario que la tabla esté registrada en un catálogo interno empleado durante la ejecución del programa. La funcionalidad de consulta empleando SQL aún está en *beta* y no tiene todas las posibilidades de consulta implementadas.

Se puede realizar una abstracción de una tabla a partir de diversas fuentes de origen empleando *TableSource*: bases de datos (por ejemplo: MySQL, HBase), ficheros (csv) o soluciones de mensajería (por ejemplo: Apache Kafka, RabbitMQ).

- **Gelly** [46]: esta librería está concebida para el procesamiento en batch de grafos. Se proporcionan métodos para crear, transformar y modificar los grafos junto a una biblioteca de algoritmos típicos utilizados en el tratamiento de grafos.

4.2.2 Lenguaje R y R Studio

R [32] es un lenguaje de programación y entorno especializado en análisis estadístico. Además de ser software libre y de disponer de multitud de herramientas propias cuenta con una amplia comunidad de usuarios que han desarrollado y puesto a disponibilidad pública multitud de paquetes y extensiones que amplían aún más las capacidades de este lenguaje y facilitan el desarrollo de experimentos y cálculos (CRAN: *Comprehensive R Archive Network*, es el repositorio de paquetes centralizado para R).

R destaca por su facilidad para integrarse con diferentes bases de datos y trabajar con diferentes formatos de fichero para la carga y lectura de datos sobre los que se necesite trabajar. Ofrece una potente herramienta de generación y personalización de gráficos y estos se pueden enriquecer aún más con muchos de los paquetes desarrollados especializados en representación y visualización de datos.

La versión empleada para la implementación en este proyecto ha sido la 3.2.3.



Figura 4.10 - Logotipos de R y R Studio

R Studio [33] es un entorno de desarrollo integrado (del inglés, *IDE*, *Integrated Development Environment*) de código abierto pensado para trabajar de manera sencilla con R. Dispone de una consola donde se pueden ejecutar comandos, un editor de código que facilita la escritura de *scripts* complejos -facilitando la identificación de errores y su ejecución directa-, una herramienta de depuración y una ventana de acceso a los gráficos generados y a la ayuda en línea. Desde esta solución gráfica se gestiona fácilmente la instalación de nuevos paquetes desde CRAN y sus dependencias. R Studio se ofrece de manera gratuita aunque existe una versión comercial para empresas con funcionalidades avanzadas y soporte.

JRI [34] es un conjunto de librerías que implementan una interfaz entre Java y R, facilitando la ejecución de sentencias en el entorno R desde una aplicación Java. Instalando el paquete *rJava* en R, ejecutando el script de configuración que se descarga junto con la instalación y enlazando las librerías de R en Java apropiadas al proyecto

(*JRI.jar*, *JRIEngine.jar*, *REngine.jar*) se puede invocar el intérprete de R desde una clase Java.

4.2.3 Keras y TensorFlow

TensorFlow [30] es una librería de código abierto originariamente desarrollada por ingenieros del equipo del *Google Brain Team* [59] buscando facilitar la investigación en aprendizaje automático y haciendo fácil y rápido el paso del prototipo desarrollado a un entorno productivo. Está disponible desde Noviembre de 2015 bajo licencia Apache 2.0 y continuamente recibe nuevas aportaciones de otros usuarios.

TensorFlow se basa en la construcción de grafos para el flujo de datos. Los nodos en el grafo representan operaciones matemáticas, entradas/salidas de datos o variables. Las aristas entre ellos representan *tensors*, relaciones de comunicación entre los nodos que no son más que matrices multidimensionales de datos que se van pasando de nodo en nodo. Está concebido para ejecutarse en múltiples dispositivos aprovechando las capacidades específicas de procesamiento de cada uno de ellos permitiendo una API única independiente del tipo de dispositivo. A los nodos se les asignan los recursos computacionales existentes y se va realizando su ejecución de forma asíncrona y paralela, a medida que los datos de los que se alimentan vayan estando disponibles.

TensorFlow no pretende ser una librería exclusiva de redes neuronales; en realidad, cualquier programa que pueda ser expresado como un grafo por donde fluya la información se puede representar con esta librería. Se ofrece por defecto una colección de módulos que se pueden reutilizar en las estructuras de redes de neuronas más comunes, pero también se pueden escribir en Python o C++ nuevas funcionalidades que se reutilicen en modelos más complejos.

Existen numerosos ejemplos y tutoriales en la propia página oficial de TensorFlow pero si no se es un experto en la materia llega a ser bastante costoso el familiarizarse con la terminología y su API. Una forma mucho más sencilla de beneficiarse del potencial de TensorFlow pero utilizando una interfaz más intuitiva es Keras, que se explica a continuación.

TensorFlow no está disponible para ser instalado en Windows, así que se ha decidido instalarlo en Docker -sobre Windows-.



Figura 4.11 - Logos de TensorFlow y Keras

Keras [31] es una librería desarrollada en Python por el ingeniero François Chollet que es capaz de ejecutar modelos de aprendizaje profundos sobre Theano o TensorFlow. La principal ventaja que ofrece sobre el uso directo de estas herramientas es que su sintaxis es mucho más sencilla e intuitiva.

Los principios sobre los que se ha desarrollado Keras son:

- **Modularidad:** cada modelo se construye a partir de un conjunto de piezas que van encajando unas con otras, tratando de minimizar las restricciones en su integración. Al final, el modelo se convierte en una secuencia en la que se van añadiendo las capas de neuronas, funciones de coste, optimizadores, funciones de activación y todo el resto de elementos que sean necesarios en el modelo concreto a implementar.
- **Minimalista:** cada modulo desarrollado en Keras es lo más reducido y simple posible, buscando facilitar la integración y la velocidad en la compilación del modelo.
- **Facilidad** para la creación de nuevos módulos: se pueden incluir nuevos componentes más avanzados como funciones o clases a los ya existentes.
- Todo está en **Python:** no hay ficheros de configuración ocultos para los modelos. Todos los módulos están escritos en Python, facilitando su depuración y su inclusión en modelos posteriores más sofisticados.

Las versiones de Python soportadas por Keras van desde la 2.7 a la 3.5.

Por defecto, tras su instalación Keras se ejecuta sobre Theano. Para seleccionar el *backend* sobre el que se ejecuta Keras es necesario editar el fichero *keras.json* que se almacena en el directorio *.keras* en la ruta del usuario modificando el valor del parámetro '*backend*' fijándolo a *tensorflow* en lugar de a *theano*.

4.2.4 Docker

Docker [35] permite la ejecución de aplicaciones dentro de *contenedores* software. Un contenedor incluye la aplicación y todas sus dependencias, pero comparte el kernel del sistema operativo con otros contenedores, ejecutándose como una aplicación aislada del resto de contenedores. Es similar a una maquina virtual pero sin necesidad de incluir las librerías y binarios del sistema operativo.



Figura 4.12 - Logotipo de Docker

La ventaja de Docker es que permite portar la aplicación a otros entornos donde se pueda ejecutar Docker: en la Figura 4.13 se muestran las diferencias de componentes necesarios a la hora ejecutar una aplicación sobre una máquina virtual o sobre Docker.

En este proyecto el contenedor se ejecuta sobre una maquina virtual Linux; junto con el instalador de Docker se incluye también *Oracle VM Virtual Box* para la gestión de máquinas virtuales. Se ejecuta el contenedor de Docker con la *imagen* de TensorFlow y se instala Keras en él.

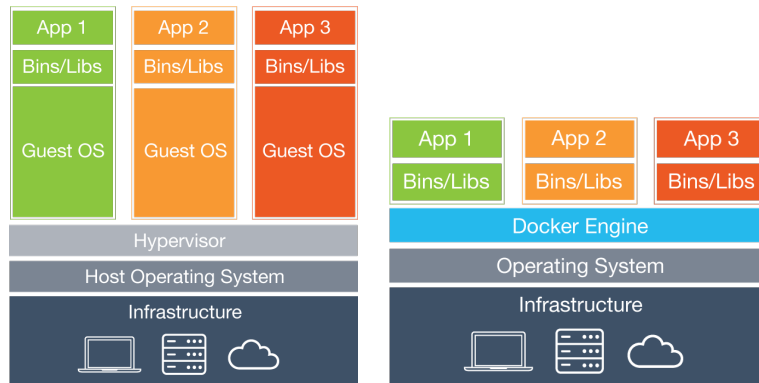


Figura 4.13 - Comparativa aplicación sobre máquina virtual y sobre Docker

Para la ejecución de Docker en local, como se hace en la fase de implementación de este proyecto, es gratuito su uso. Existe una modalidad en *cloud* denominada *Containers-as-a-Service* (CaaS) que sí que exige pago por suscripción.

4.2.5 Python

Python [42] es un lenguaje de programación creado a principios de la década de los 90 por Guido van Rossum que buscaba desarrollar un lenguaje con una sintaxis que favoreciera el desarrollo de código fácilmente legible.

Se trata de un lenguaje interpretado, multiparadigma (soporta Orientación a Objetos y programación imperativa) y multiplataforma. El tipado de variables se hace de manera dinámica.

Entre las versiones 2.x y 3.x del lenguaje existen importantes diferencias que las hacen ser incompatibles en muchos aspectos y exigen la reescritura de buena parte del código antiguo.

Python es un lenguaje que se ha vuelto muy popular debido a su facilidad de aprendizaje –de hecho, es uno de los lenguajes más utilizados para iniciarse en la programación–, goza de una amplia comunidad de usuarios y de un uso extendido en el desarrollo web.

El uso de Python en este proyecto está limitado a la fase de construcción y modelado de la red de neuronas empleando además las librerías de Keras y TensorFlow.



Figura 4.14 - Logotipo de Python

4.2.6 Java y Eclipse

Java [43] es un lenguaje de programación Orientado a Objetos (*OOP*, *Object Oriented Programming*), desarrollado por Sun Microsystems en los años 90. Mucha de su sintaxis deriva de C y C++, pero se simplifica el modelo de datos y se eliminan herramientas de bajo nivel que suelen inducir a errores -como la manipulación directa de punteros o de memoria-.

Típicamente, las aplicaciones en Java se escriben una sola vez y se pueden ejecutar en cualquier plataforma donde se ejecute la JVM (*Java Virtual Machine*, Máquina Virtual de Java). Las clases de código Java se compilan en un *bytecode* y éste es interpretado a código nativo para la ejecución en una plataforma concreta gracias a la máquina virtual, que transforma el *bytecode* a instrucciones específicas del entorno de ejecución facilitando su portabilidad entre diferentes plataformas.

Es un lenguaje de tipado estático. Java dispone de una gran colección de bibliotecas y de un API diferenciada según el tipo de dispositivo al que se oriente el uso del programa a desarrollar. En la JDK (*Java Development Kit*, distribuida por Oracle, ahora propietaria de la originaria Sun Microsystems) se incluye el compilador y el entorno de ejecución (*JRE*, *Java Runtime Environment*), el depurador y el generador de documentación, entre otros.

Actualmente Java se distribuye bajo Licencia Pública General de GNU, haciendo su uso gratuito para la fase de implementación de este proyecto. Para la ejecución de las aplicaciones Java utilizando Apache Flink se ha utilizado la versión de JDK SE (*JDK Standard Edition*) 8.



Figura 4.15 - Logotipos de Java y Eclipse

Como entorno de desarrollo de las aplicaciones Java se ha empleado Eclipse Mars, que es la distribución de Eclipse que más se ajusta a las necesidades de desarrollo con Apache Flink (facilitando, por ejemplo, la integración con Maven). Eclipse proporciona una rica interfaz de desarrollo desde donde escribir, compilar, ejecutar y depurar programas desarrollados en lenguaje Java. Está liberado bajo licencia de software libre así que no supone ningún gasto añadido a la elaboración del proyecto, en que se ha utilizado la versión de Eclipse Mars.1 Release (4.5.1).

5. Implementación

5.1 Introducción

En este apartado se van a desglosar las diferentes partes que se han ido desarrollando para resolver el problema planteado, partiendo de los ficheros de contratos iniciales hasta conseguir los resultados de las predicciones finales.

Se irán detallando cada uno de los módulos desarrollados en el orden cronológico que han sido implementados para entender mejor la necesidad que pretende satisfacer cada uno de ellos.

5.2 Análisis inicial de los ficheros

El punto de partida son los ficheros con la información de contratación de la administración pública de EE. UU. entre el año fiscal del 2000 y el año fiscal 2015.

En [23] se detalla el formato y los campos en orden alfabético. Son ficheros de texto plano con extensión `.csv` y con una línea de cabecera. Cada uno de los campos está separado del siguiente por una coma `,`, y cada campo aparece entrecomillado por comillas dobles (`"campo1"`, `"campo2"`). En la Figura 5.1 se incluye una imagen donde se muestra un extracto de ejemplo de uno de los ficheros de contratación donde se puede apreciar su formato.

Se incluyen 250 campos -características- en cada una de las líneas de un fichero. Cada una de estas líneas hace referencia a los detalles de adjudicación de un contrato. Recordando que el objetivo es poder predecir cuándo se va a adjudicar un contrato en un sector determinado, se ha realizado una clasificación de estas características para ver cuáles de ellas podrían utilizarse como indicadores de una adjudicación.

```

1 "unique_transaction_id","transaction_status","dollarsobligated","baseandexercisedoptionsvalue","base
2 "7eb45398-0d11-434a-9059-88945d6456b9","Active","7963720.2200","7963720.2200","7963720.2200","9700:
3 "576bb0a7-201f-4056-a76f-b00aa2916fbc","Active","36068.5600","36068.5600","36068.5600","9700: Depart
4 "e6a49552-9b7e-4394-9d12-6b8db1aba6af","Active","223641.5300","223641.5300","223641.5300","9700: Dep
5 "e3f05860-cf79-4b28-afcf-2d9396312de7","Active","3748721.8900","3748721.8900","3748721.8900","9700:
6 "9343fd1f-434d-463c-a55f-84c9fbc52cff","Active","13945.5000","13945.5000","13945.5000","9700: Depart
7 "4f76a1a5-ad2f-482e-80ef-9a7100b8a382","Active","51.4400","51.4400","51.4400","9700: Department of D
8 "9163ed4a-3f6b-4155-acf4-31e84a585254","Active","13945.5000","13945.5000","13945.5000","9700: Depart
9 "4939a508-8df2-4962-b99c-7469ad79b471","Active","108919.1400","108919.1400","108919.1400","9700: Dep
10 "bc30d36e-9ee0-4041-afeb-939be23c16b3","Active","2058.8300","2058.8300","2058.8300","9700: Departmen
11 "c937b848-fa11-4034-a4ba-30528d260eb8","Active","132531.1900","132531.1900","132531.1900","9700: Dep

```

Figura 5.1 - Extracto de un fichero de texto con información de los contratos

Se han diferenciado tres grandes conjuntos de categorías:

- **Indicadores administrativos:** *a76action*, *davisbaconact*, *idvpiid*, etc. Se incluyen aquí aquellos campos que hacen referencia a identificadores internos de la administración y a marcadores que informan sobre si el contrato se ajusta a determinadas leyes o directrices de adjudicación.
- **Indicadores sobre el adjudicatario:** *isforeigngovernment*, *haobflag*, *prime_awardee_executive1*, *prime_awardee_executive1_compensation*, etc. Se

incluyen en esta categoría todos aquellos campos que hacen referencia a cualidades del prestador del servicio al que ha sido adjudicado el contrato. Hay marcadores que indican el origen del proveedor, el nombre e ingresos de sus ejecutivos principales, datos fiscales, etc. Como estos datos sólo pueden saberse una vez se ha adjudicado el contrato, carece de sentido utilizarlos en la predicción de cuándo se producirá una próxima adjudicación.

- **Indicadores de las características del contrato:** aquellos que recogen alguna propiedad del contrato y que se conoce en el momento de su licitación. Dado que se parte de contratos ya adjudicados y no de anuncios de licitación, estas propiedades serán las más útiles a la hora de identificar un contrato y poder emplearse para realizar un pronóstico.

Las 19 características que a priori parecen más interesantes son:

- ***agencyid*** [Formato: Cadena de texto libre, 500 caracteres máximo]: identificador de la Agencia u Oficina que es responsable de la adjudicación del contrato.
- ***dollarsobligated*** [Formato: Numérico. 8 dígitos máximo. Puede ser positivo o negativo. Se emplea el punto ‘.’ como separador decimal]: el importe neto de dólares de la transacción. Si se trata de una modificación o rescisión este importe puede ser menor o igual a 0.
- ***claimantprogramcode*** [Formato: Cadena de texto libre, 500 caracteres máximo]: código del programa solicitante que identifica el grupo de suministros, construcción u otro tipo de servicios demandados.
- ***contractactiontype*** [Formato: Cadena de texto libre, 500 caracteres máximo]: tipo de concesión. Estos tipos incluyen: *Purchase Orders (PO)*, *Delivery Orders (DO)*, *BPA Calls* y *Definitive Contracts*.
- ***contractingofficeagencyid*** [Formato: Cadena de texto libre, 500 caracteres máximo]: la definición que se da para este campo es la misma que en el *agencyid*.
- ***contractingofficeid*** [Formato: Cadena de texto libre, 500 caracteres máximo]: código de la oficina de contratación.
- ***descriptionofcontractrequirement*** [Formato: Cadena de texto libre, 5.000 caracteres máximo]: una breve descripción de los bienes o servicios demandados.
- ***fiscal_year*** [Formato: Cadena de texto libre, 4 caracteres máximo]: año fiscal.
- ***fundingrequestingagencyid*** [Formato: Cadena de texto libre, 500 caracteres máximo]: código que identifica la agencia que aporta la mayoría de los fondos para el pago del contrato.

- ***multipleorsingleawardidc*** [Formato: Cadena de texto libre, 500 caracteres máximo]: indica si el contrato pertenece a una asignación simple o a un grupo de contratos adjudicados conjuntamente.
- ***multiyearcontract*** [Formato: Cadena de texto libre, 500 caracteres máximo]: indica si el contrato tiene una duración de entre 1 y 5 años.
- ***principalnaicscode*** [Formato: Cadena de texto libre, 50 caracteres máximo]: indica el código NAICS (*North American Industry Classification System*, usado para identificar el sector económico al que se dedica una empresa en las economías de México, Canadá, y Estados Unidos). A priori es un campo que está relacionado con el adjudicatario del contrato, pero que también lo está con el objeto del contrato y los bienes o servicios demandados.
- ***productorservicecode*** [Formato: Cadena de texto libre, 500 caracteres máximo]: código que mejor identifica el producto o servicio demandado.
- ***renewaldate*** [Formato: Cadena de texto libre, 12 caracteres máximo. El formato de la fecha es *mm/dd/aaaa* (mes, día, año)]: no se ofrece una descripción oficial, pero se puede entender como la fecha de renovación, prórroga o extensión de un contrato.
- ***signeddate*** [Formato: Cadena de texto libre, 12 caracteres máximo. El formato de la fecha es *mm/dd/aaaa* (mes, día, año)]: fecha de adjudicación de un contrato.
- ***reasonformodification*** [Formato: Cadena de texto libre, 500 caracteres máximo]: en caso de que el registro se trate de una modificación en este campo se incluye el código de identifica el tipo de modificación.
- ***typeofcontractpricing*** [Formato: Cadena de texto libre, 500 caracteres máximo]: código para el tipo de contrato en relación a si se trata de un precio fijo o si depende del coste en que se incurra.
- ***ultimatecompletiondate*** [Formato: Cadena de texto libre, 12 caracteres máximo. El formato de la fecha es *mm/dd/aaaa* (mes, día, año)]: fecha estimada o planificada del final del contrato, incluyendo la posible ejecución de todas las ampliaciones provistas en el contrato.
- ***modnumber*** [Formato: Cadena de texto libre, 50 caracteres máximo]: código que identifica unívocamente una modificación para un contrato.

5.3 Preselección de características

Una vez identificadas las 19 características de las 250 totales que podrían ser de utilidad como identificadores de un contrato lo siguiente que se va a hacer es

transformar los ficheros originales en otros más pequeños que sólo incluyan las características preseleccionadas. Se generará un fichero de salida por cada año fiscal y se estudiará cada uno de los campos.

Para la preselección de características se va a utilizar la clase *Preseleccion.java*. En esta implementación se utiliza Flink para seleccionar aquellas columnas que se desea conservar de los ficheros origen.

Por cada uno de los ficheros se hace:

- Se fija como fuente (*source*) el fichero *.csv* de contratos del año a tratar. (*readCsvFile*).
- Puesto que la primera línea es de cabecera con el nombre de los campos, se ignora (*ignoreFirstLine*).
- En caso de que alguna línea no se ajuste al formato general, se ignorará también (*ignoreInvalidLines*).

Hay que indicar para cada columna si se quiere incluir en el *dataset* resultante de la lectura del fichero *.csv* o no. Para ello existe el método *includeFields* que recibe como parámetro una cadena de texto de 0s y 1s indicando por cada posición si se desea conservar ese campo en el *dataset* generado a partir del fichero. En este caso va a ser necesaria una cadena de 250 dígitos con 19 posiciones, las correspondientes a las características que se quieren conservar, fijadas a 1. Para evitar lo tedioso que sería construir esa cadena de forma manual y además así poder generalizar el procedimiento en caso de que en un futuro se incluyeran o eliminaran campos, se construye esa cadena de manera dinámica. Se han creado para ellos dos ficheros:

- Un fichero que contiene la primera línea de cabecera de cualquier año fiscal. Esa línea contiene los nombres en orden de los campos según se ha construido el fichero.
- Un fichero de texto plano que contiene en cada línea el nombre de una característica a incluir en el fichero de salida.

Se carga el contenido de ambos ficheros en memoria y si el campo en el fichero de cabecera coincide con alguna de las características incluidas en el otro fichero se marca un 1 en la cadena; 0 en caso contrario.

Para la carga de estos ficheros de cabecera y características se utiliza también Flink, aunque son pequeños y podrían leerse de otro modo.

Una vez que se ha construido la cadena con los campos a incluir se le puede indicar a la fuente que los campos vienen parseados con comillas dobles, con el método *parseQuotedStrings*. Esta opción ha presentado dos problemas y por ello se ha descartado:

- Si en algún momento alguna línea no cumple con todos los campos parseados correctamente la aplicación genera una excepción.
- Si alguno de los campos a leer se define como numérico pero está parseado, Flink primero intenta cargarlo como numérico por lo que lanza una excepción `NUMERIC_VALUE_FORMAT_ERROR`.

Además nos encontramos con otro problema añadido: el método de cargar el fichero elegido es un *readCsvFile* que por defecto interpreta que cada campo está separado por una coma ','. Lo que ocurre es que algunos de los campos contiene comas en su contenido (por ejemplo, la agencia "3600: VETERANS AFFAIRS, DEPARTMENT OF" lo que provoca un fallo al tratar de cargar la línea del fichero.

Para simplificar la carga y disminuir el número de líneas consideradas como inválidas se decide no utilizar el método *parseQuotedStrings* y utilizar en su lugar un delimitador de campo más adecuado a este caso: ";", indicado en el método *fieldDelimiter("\"\", \"\")*. Con esto se eliminan las dobles comillas y la coma que separan cada campo, aunque quedarían por eliminar las dobles comillas del comienzo del primer campo y del final del último campo. Para esta última acción se aplica una transformación *map* que elimine las dobles comillas y que sustituya las posibles tabulaciones por espacios en blanco.

Una vez transformado el *dataset* se vuelcan las tuplas en orden a un fichero de texto plano separado por tabuladores.

Para la realización del proceso se puede indicar a Flink el índice de paralelismo que se puede emplear para realizar las operaciones sobre cada fichero. El equipo empleado dispone de 4 procesadores lógicos, por lo que se podría indicar este número como índice de paralelización y que cada fichero se dividiera en 4 particiones y que sobre cada una de ellas se realizarán las transformaciones necesarias hasta generar el fichero de salida. A la hora de volcar el *dataset* a un fichero de salida se generan tantos ficheros de salida como particiones se hayan empleado; para evitar que por cada fichero de entrada del año fiscal se generen 4 ficheros de salida distintos, se fija el grado de paralelización en la escritura del fichero a 1. Así que se disponen de dos índices de paralelización diferenciados: el global para el tratamiento de particiones y el particular fijado para el volcado a fichero.

El fichero de salida también va a incluir una línea de cabecera con el nombre de cada una de las características seleccionadas que identifique la columna de datos. Para asegurar que la cabecera se incluya en primera posición y no en otra cuando las particiones se unan en una sola para la escritura, se incluye una operación de ordenación que asegure que la primera línea es la de cabecera.

El resultado de este procesado de los 15 ficheros se recoge en la Tabla 5.1

Nombre del fichero	Líneas totales	Líneas procesadas	Inválidas
2000_All_Contracts_Full_20160115.csv	594.599	594.598	1
2001_All_Contracts_Full_20160115.csv	642.063	642.063	0
2002_All_Contracts_Full_20160115.csv	830.600	830.600	0
2003_All_Contracts_Full_20160115.csv	1.183.843	1.183.842	1
2004_All_Contracts_Full_20160115.csv	2.002.015	2.001.916	99
2005_All_Contracts_Full_20160115.csv	2.923.305	2.923.046	259
2006_All_Contracts_Full_20160115.csv	3.797.163	3.797.061	102
2007_All_Contracts_Full_20160115.csv	4.111.343	4.111.302	41
2008_All_Contracts_Full_20160115.csv	4.504.817	4.504.779	38
2009_All_Contracts_Full_20160115.csv	3.496.847	3.496.772	75
2010_All_Contracts_Full_20160115.csv	3.538.783	3.538.733	50
2011_All_Contracts_Full_20160115.csv	3.395.661	3.395.632	29
2012_All_Contracts_Full_20160115.csv	3.116.221	3.116.197	24
2013_All_Contracts_Full_20160115.csv	2.504.847	2.504.821	26
2014_All_Contracts_Full_20160115.csv	2.512.864	2.512.847	17
2015_All_Contracts_Full_20160115.csv	3.639.204	3.639.012	192
Total	42.794.175	42.793.221	954

Tabla 5.1 - Número de líneas por cada fichero de año fiscal

El cómputo total de líneas incluye la línea de cabecera tanto en los ficheros de entrada como en los ficheros resultados de la preselección. Aquellas líneas que no han sido incluidas en el fichero de salida es porque se han considerado inválidas al no ajustarse al formato general especificado para el fichero.

5.4 Estudio de las características preseleccionadas

De cada una de las características preseleccionadas se quiere conocer cuál es el dominio en el que toma valores y su número de repetición a lo largo de los años. Para hacer este tratamiento más estadístico se ha creado un script utilizando R como lenguaje y R Studio como entorno de trabajo. Se ha cargado en memoria cada uno de los ficheros creados en el paso anterior y, por cada una de las características se ha generado una tabla de frecuencias. Al final del proceso se unirán todas esas tablas de frecuencias en una única para sacar el cómputo total de todos los años.

Teniendo en cuenta que en total, desde 2000 a 2015 se tienen más de 44 millones de registros sobre contratos y consultando el detalle de cada característica se observa que:

- **agencyid:** todos los registros lo tienen informado. Hay 290 valores distintos. El más repetido, con más de 19 millones de registros, es la agencia "9700: DEPT OF DEFENSE".

Los literales del *agencyid* tienen un código numérico seguido de dos puntos ':' y una cadena de texto a continuación. Se ha considerado que dos literales de *agencyid* son el mismo (y por tanto se codificarán con el mismo valor numérico para la predicción) si el código numérico coincide. Así, por ejemplo, los pares de *agencyid*: "4745: OFFICE OF GOVERNMENT WIDE POLICY" y "4745: OFFICE OF THE CHIEF ACQUISITION OFFICER", "12B0: USDA, ASSISTANT SECRETARY FOR ADMINISTRATION" y "12B0: USDA, ASSISTANT SECRETARY FOR DEPARTMENTAL MANAGEMENT", "7528: AGENCY FOR HEALTH CARE POLICY AND RESEARCH" y "7528: AGENCY FOR HEALTHCARE RESEARCH AND QUALITY", "1406: OFFICE OF POLICY, BUDGET AND ADMINISTRATION" y "1406: OFFICE OF POLICY, MANAGEMENT, AND BUDGET", "7523: CENTERS FOR DISEASE CONTROL" y "7523: CENTERS FOR DISEASE CONTROL AND PREVENTION" se han considerado equivalentes.

Hay otras *agencyid* que aunque su literal textual es idéntico, su codificación numérica no lo es, y se ha decidido mantenerlos como *agencyid* diferentes. Ejemplos: "7004: OFFICE OF THE INSPECTOR GENERAL" y "1404: OFFICE OF THE INSPECTOR GENERAL", o "4730: FEDERAL ACQUISITION SERVICE" y "4732: FEDERAL ACQUISITION SERVICE".

Esta característica se considera que puede resultar de utilidad a la hora de identificar un contrato, pues no está excesiva fragmentada y aporta información sobre el órgano que solicita la licitación.

- **dollarsobligated:** esta característica toma valor en el rango numérico continuo en el intervalo [-344.733.702.433, 344.733.911.304] dólares. La media del importe de los contratos es de 155.989,90\$.

La distribución del importe de los contratos se puede apreciar en diagrama de caja de la Figura 5.2.

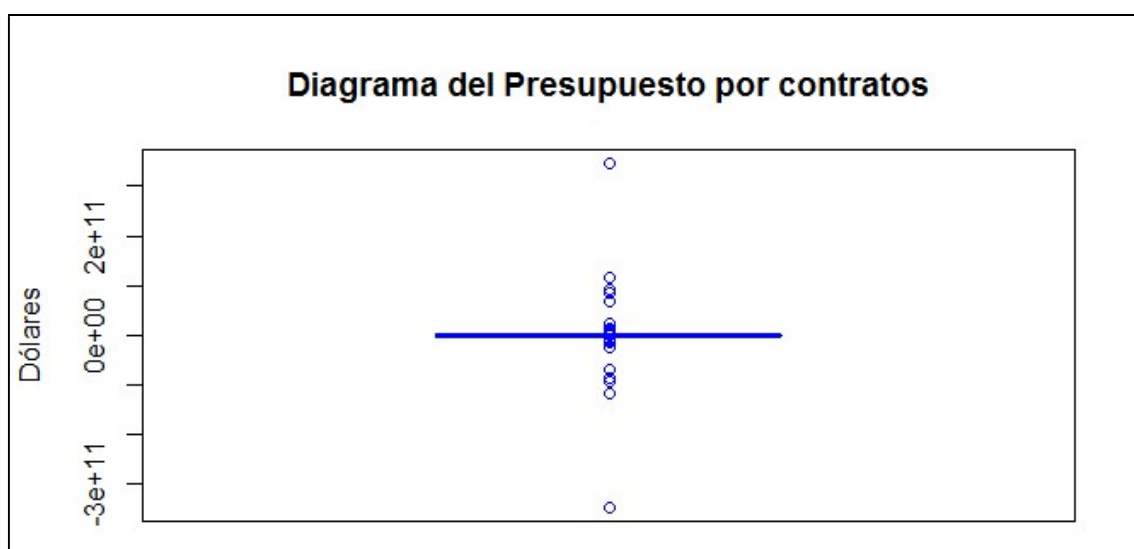


Figura 5.2 - Distribución de los contratos según su presupuesto

Se aprecia que el mínimo y el máximo del rango indicados constituyen los dos valores más extremos (en inglés, *outliers*) de la muestra. Si se eliminan esos dos, el rango en el que se toman valores es: [-116.092.075.155, 116.091.716.928]. En el gráfico de la Figura 3.5 aún se puede apreciar que son valores extremos aunque no tanto como los anteriores. En cualquier caso lo interesante de estos valores es que es bastante probable que constituyan *ruido* en los datos: no se trata de contratos reales o de cancelaciones -los de importe negativo- de cientos de miles de millones de dólares: se trata de cifras que han sido mal grabadas, de ahí que aparezca la cantidad en positivo y más tarde una cantidad en negativo que sirve para cancelar la anterior.

Por ejemplo, para el caso de los primeros valores extremos identificados se puede ver fácilmente en la Figura 5.3 de la web de UsaSpending.gov que ese ha sido el caso.

MISCELLANEOUS FOREIGN CONTRACTORS	W912GB09C0090	\$344,733,911,304	9/26/2013	Contracts	Department of Defense
MISCELLANEOUS FOREIGN CONTRACTORS	W912GB09C0090	\$215,383	9/22/2010	Contracts	Department of Defense
MISCELLANEOUS FOREIGN CONTRACTORS	W912GB09C0090	\$18,698	11/18/2013	Contracts	Department of Defense
MISCELLANEOUS FOREIGN CONTRACTORS	W912GB09C0090	(\$2,247,292)	11/7/2012	Contracts	Department of Defense
MISCELLANEOUS FOREIGN CONTRACTORS	W912GB09C0090	(\$344,733,702,433)	9/27/2013	Contracts	Department of Defense

Contract Description

ADMINISTRATIVE MODIFICATION TO CORRECT TYPOGRAPHICAL ERROR ON PRIOR MODIFICATION. THIS CONTRACT ACTION IS GOVERNED BY THE ABG75 PROCESS AND NOT THE FAR.

Figura 5.3 - Captura donde se muestran los contratos con importes más altos y más bajos. Se incluye el motivo de registro del contrato menor

Para facilitar la predicción es necesario dividir en subrangos el intervalo de valores y asignar a cada uno de ellos un identificador único numérico.

La distribución escogida para su *discretización* según los intervalos de cantidad dineraria se muestra en la Tabla 5.2, indicándose el presupuesto en dólares, el número de registros que pertenecen a ese intervalo en el periodo 2000/15 y el número entero asignado a ese intervalo para su codificación posterior.

- ***claimantprogramcode***: aparecen 33 valores distintos. El más repetido es ‘.’, que puede considerarse que el campo no está informado, y se repite para más de 23 millones de contratos. Dado el elevado número de registros que no lo tienen informado no va a considerarse como una característica reseñable.
- ***contractactiontype***: aparecen 102 valores. De ellos, sólo 8 valores se aplican a más de 100.000 registros, siendo "PO Purchase Order", con más de 11 millones y "DO Delivery Order", con más 23 millones, los más repetidos.
- ***contractingofficeagencyid***: tiene una distribución muy parecida al *agencyid*. Hay 323 valores distintos y cada valor se compone de un número separado por dos puntos de una cadena de texto. Algunos tienen el mismo código numérico pero un literal diferenciado, por lo que se han considerado similares (si no se hubiera hecho, habría un total de 333).

Rango (en \$)	N.º de Contratos	Codificación
Negativos (<0)	2.470.912	1
Iguales a 0	5.372.703	2
(0,3.000]	12.261.680	3
(3.000, 10.000]	7.948.435	4
(10.000, 25.000]	4.824.300	5
(25.000, 50.000]	3.177.140	6
(50.000, 75.000]	1.426.413	7
(75.000, 100.000]	1.004.122	8
(100.000, 150.000]	949.486	9
(150.000, 250.000]	964.292	10
(250.000, 500.000]	983.524	11
(500.000, 750.000]	383.902	12
(750.000, 1.000.000]	228.620	13
(1.000.000, 2.000.000]	358.986	14
>2.000.000	437.870	15
No numéricos (NaN)	820	

Tabla 5.2 – Valores discretos para los rangos de presupuesto

- ***contractingofficeid***: hay más de 7.600 valores diferentes, por lo que estando tan fragmentado se prefiere no incluirlo en las características definitivas.
- ***descriptionofcontractrequirement***: se trata de una característica de texto libre, por lo que es difícil categorizarla. Se descarta utilizarla como predictor.
- ***fiscal_year***: se busca precisar la fecha de adjudicación por mes y año. Al comenzar en octubre, dentro un mismo año fiscal hay meses de dos años naturales diferentes, por lo que se declina utilizar esta característica y se opta por obtener el mes y año del campo *signeddate*.
- ***fundingrequestingagencyid***: aparecen casi 900 valores distintos y más de 9 millones de registros no la tienen informada. Se descarta su utilización como predictor.
- ***multipleorsingleawardidc***: se informan 7 valores distintos, pero más de 41 millones de registros no lo tienen informado. Se descarta su empleo como característica definitiva.
- ***multiyearcontract***: hay 11 valores distintos. Así 7 millones de registros no lo tienen informado y 35 millones tienen informado el mismo valor. Se descarta su uso como característica definitiva.

- ***principalnaicscode***: hay 1.637 valores distintos (habría 1.664 de no ser por haber considerado repeticiones si coincide su código numérico como en el caso del *agencyid*). Tres millones de registros no lo tienen informado, pero se ha comprobado que en su mayoría son modificaciones sobre contratos y no adjudicaciones nuevas.

A pesar de ello se considera un buen indicador del tipo de proveedor que puede estar interesado en un contrato determinado.

- ***productorservicecode***: esta característica hace referencia al servicio demandado en vez de al proveedor. Hay 200.000 registros que no lo tienen informado pero son más de 5.000 valores distintos y hay repeticiones de código con distinta descripción. Se declina utilizar esta característica por su alta fragmentación en favor del código NAICS.

- ***renewaldate***: toma valores entre el 1 de enero de 1900 y el 1 de octubre de 2017. Hay bastantes registros que no lo tienen informado y la fecha del 01/01/1900 parece una por defecto, puesto que la siguiente fecha mayor es el 1 de junio de 1999. No está clara la utilidad de este campo ni su relación con el resto y su descripción tampoco resuelve la duda; se descarta su uso como característica.

- ***signeddate***: de esta fecha se extraerá el mes y el año en el que se adjudica el contrato. Esos dos valores marcan el orden en la serie temporal, así que serán características definitivas.

- ***reasonformodification***: 30 valores distintos pero 30 millones de registros no lo tiene informado. Parece que sólo es de aplicación a oportunidades en curso, contratos ya adjudicados.

- ***typeofcontractpricing***: aparecen 143 valores distintos. 28 millones de registros aplican el mismo valor.

- ***ultimatecompletiondate***: aparecen valores en el rango entre 1931 y 2108. No está clara su utilidad, así que se descarta su uso como característica determinante.

- ***modnumber***: aparecen más de 77.474 valores distintos. No obstante, el valor '0' se repite para más de 30 millones de registros. Ocurre que sobre un contrato asignado se pueden realizar correcciones, -porque alguno de los datos del registro anterior reportado estuvieran mal-, ejecutar ampliaciones que estuvieran contempladas en el acuerdo original o rescindir el contrato o partidas parciales concretas. Estas modificaciones se incluyen en el fichero como registros nuevos: mantienen el identificador del contrato en los campos *piid* o *idvpiid* según corresponda, y se incluye un ordinal o texto indicativo -no hay un criterio fijo definido- para conocer cuál es la siguiente modificación. En los contratos de IDV (*Indefinite Delivery Vehicle*) el secuencial se sigue en el campo *idvmodificationnumber*, en el resto en este *modnumber*. Todos los contratos, independientemente de su tipo, tienen como *modnumber* el 0 cuando se trata de la

adjudicación, así que aunque esta característica no se utilice para predicción sí que va a servir para poder distinguir qué registros se corresponden con nuevos contratos y qué registros se corresponden con cambios sobre contratos ya asignados.

5.5 División contratos

Según lo deducido del estudio de las características escogidas, el *modnumber*, el *reasonformodification* y el signo y valor de las cantidades de *dollarsobligated* pueden ayudar a identificar si un registro se trata de una nueva adjudicación o de una modificación sobre un contrato ya existente. No todas las modificaciones tienen por qué incluir valor en *reasonformodification*, y la cantidad expresada en *dollarsobligated* puede ser positiva y tratarse aún así de una modificación (por ejemplo, en ampliaciones, contratos por consumo o ejecución de oportunidades contempladas en la adjudicación). Para poder distinguir si una línea corresponde a un nuevo contrato o a una modificación con plenas garantías hay que fijarse en el campo *modnumber*. Éste toma cualquier valor distinto de 0 en caso de tratarse de una modificación posterior a la adjudicación y de 0 para nuevas oportunidades, así que en base a esta característica se van a dividir los registros en dos partes: nuevas oportunidades y modificaciones. El objetivo del proyecto es predecir nuevas concesiones que puedan ser de interés para un proveedor interesado en participar en algún concurso relacionado con un contrato de su sector, por lo que difícilmente las modificaciones sobre contratos ya en curso pueden aportar información de utilidad para esa predicción; por otro lado, el comportamiento de las modificaciones de un contrato, el conocer cómo evolucionan los importes en cantidad y periodicidad durante el tiempo que ese contrato está vigente, se queda fuera del alcance de este trabajo.

La clase *Division.java* se encarga de la separación de los ficheros entre aquellos que registren las nuevas oportunidades de contratación y aquellos ficheros con modificaciones de contratos. Para poder estudiar cómo evolucionan ambos valores a lo largo del tiempo, tanto en número de registros como en presupuesto asignado, paralelamente a esta separación se van a generar dos agregados por cada mes y año, reportando por cada uno el número de registros que caen dentro de esa categoría y el importe total de presupuesto.

Se utilizan como ficheros de entrada los resultantes de la ejecución del proceso de Preselección, que ya sólo contienen las 19 características provisionales. Sobre el *dataset* formado con los registros de cada año fiscal se aplican dos filtros (método *filter*), implementados en las clases *NuevasOportunidadesFilter.java* y *ModificacionesFilter.java*. Un filtro devuelve verdadero o falso por cada tupla del *dataset* de entrada, de tal manera que el *dataset* de salida sólo contiene aquellas tuplas que han satisfecho la condición del filtro. Así, tras aplicar al conjunto de datos de las tuplas de los registros del año fiscal que se esté tratando el filtro de *NuevasOportunidadesFilter*, se obtiene como resultado el *dataset* con las tuplas cuyo campo *modnumber* es igual a 0. En el caso del *ModificacionesFilter* es justo al

contrario: los registros posteriores a la adjudicación del contrato –correspondientes a cambios de éste- nunca llevan un 0 informado en ese campo.

En la Tabla 5.3 se muestra el resultado de la separación de las líneas de los ficheros entre aquellas correspondientes a nuevas oportunidades y aquellas que pertenecen a modificaciones. Se ha desglosado el resultado por cada año fiscal del período 2000-2015.

Año Fiscal	Nuevas oportunidades	Modificaciones
2000	396.817	197.780
2001	433.791	208.271
2002	557.960	272.639
2003	847.260	336.581
2004	1.523.208	478.707
2005	2.283.522	639.523
2006	3.101.335	695.725
2007	3.270.628	840.673
2008	3.584.160	920.618
2009	2.491.254	1.005.517
2010	2.431.866	1.106.866
2011	2.201.178	1.194.453
2012	1.958.854	1.157.342
2013	1.395.347	1.109.473
2014	1.441.304	1.071.542
2015	2.551.497	1.087.514
Total	30.469.981	12.323.224

Tabla 5.3 - Número de nuevas contrataciones y de modificaciones

Una vez separados en dos *datasets* diferenciados los nuevos contratos y las modificaciones de contratos se construyen los agregados. Para ello, a cada *dataset* se le aplica una función *map* encargada de transformar la entrada de tuplas de 19 características en tuplas de 4, que son: un contador -a 1 para todas las tuplas-, el mes y el año -obtenidos del campo *signeddate*-, y el *dollarsobligated*.

Sobre cada uno de estos *datasets* se ejecuta otra transformación: para hacer la agregación se van a agrupar los registros por mes y año, sumando el contador y el presupuesto. Así, el *dataset* de salida contiene la información agrupada de la forma en que se buscaba.

5.6 Análisis de la evolución de la contratación

A partir de la información generada en las agregaciones de nuevas oportunidades y modificaciones obtenidas en el paso anterior, utilizando de nuevo R, se han elaborado

los siguientes gráficos para comprobar cómo ha sido la evolución de la contratación en los años fiscales de los que se disponen registros.

Para las nuevas oportunidades, en la Figura 5.4 se muestra el gráfico que refleja la evolución del número de veces que se ha producido una nueva adjudicación, mostrando el número de contratos (en miles), y en la Figura 5.5 otro gráfico por dinero invertido mostrándolo en dólares de presupuesto (en millones), agrupados por mes y año.

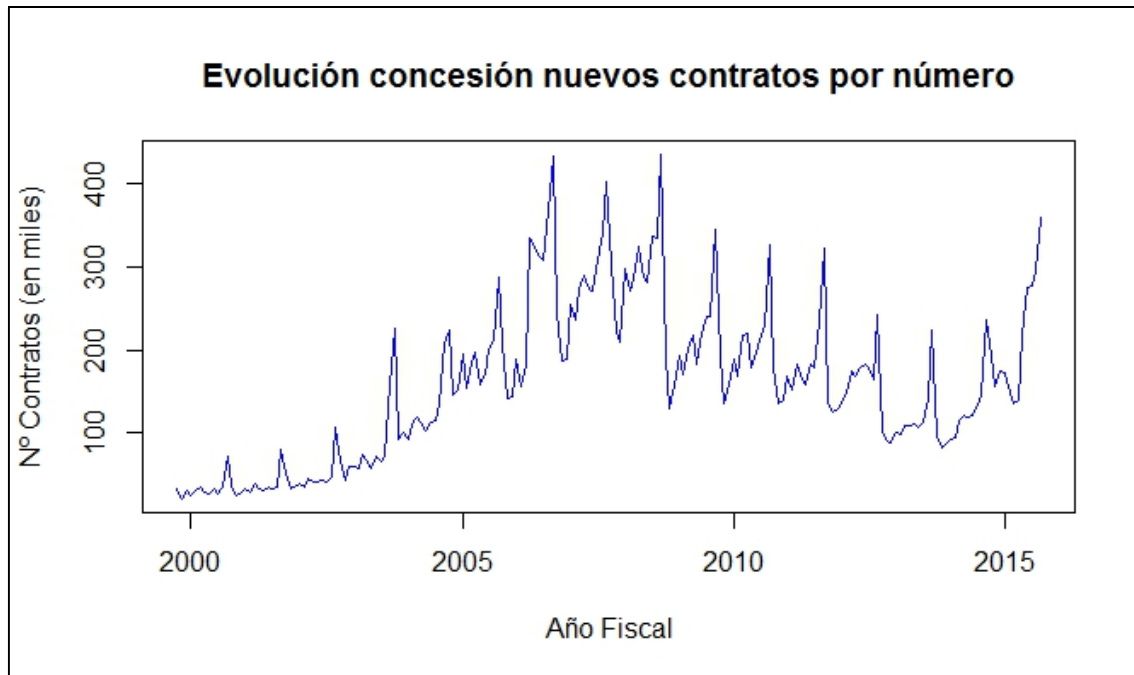


Figura 5.4 - Evolución del número de nuevos contratos (en miles) agrupados por mes para el período 2000/15

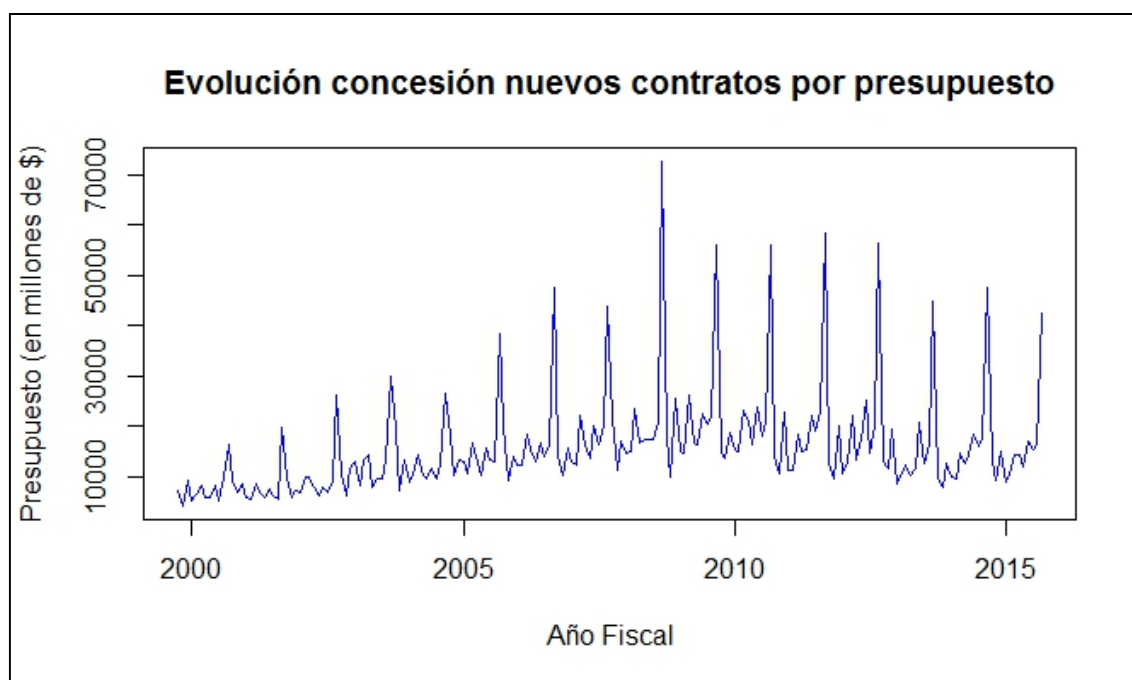


Figura 5.5 - Evolución del presupuesto destinado a nuevos contratos (en millones de \$) agrupados por mes para el período 2000/15

Los datos de los años fiscales correspondientes al período 2000-2007 durante el cual no se obligaba a reportar información, al no existir ni la página de descarga de datos ni la ley que obliga a ello, pueden no estar aún completos. De hecho, la propia web de UsaSpending.gov en sus consultas impone como límite en sus criterios de búsqueda que sea a partir del año fiscal 2008. Como desde ese año hay datos suficientes para poder elaborar este trabajo se van a descartar los ficheros anteriores a ese período, ignorando los ficheros del 2000 al 2007.

Así pues, observando en detalle el período 2008-2015 y comparándolo además con la evolución en el número de modificaciones se tiene el gráfico de la Figura 5.6. La comparación en la evolución del presupuesto destinado a nuevas contrataciones frente al destinado a modificaciones sobre contratos en vuelo se representa en el gráfico de la Figura 5.7.

En la Figura 5.6 se aprecia a simple vista que en el período 2008-2012 se da un comportamiento de ciclo estacionario, con repuntes en la contratación en los meses de abril y junio y con una fuerte subida en septiembre (último mes del año fiscal, donde podría ser que se esté consumiendo todo el presupuesto pendiente de ese ejercicio antes de empezar a aplicarse el nuevo presupuesto para el año fiscal que está a punto de comenzar). No se aprecia ninguna tendencia clara. A partir del 2012 los ciclos son menos marcados, aunque se mantiene el pico del mes de septiembre para todos los años.

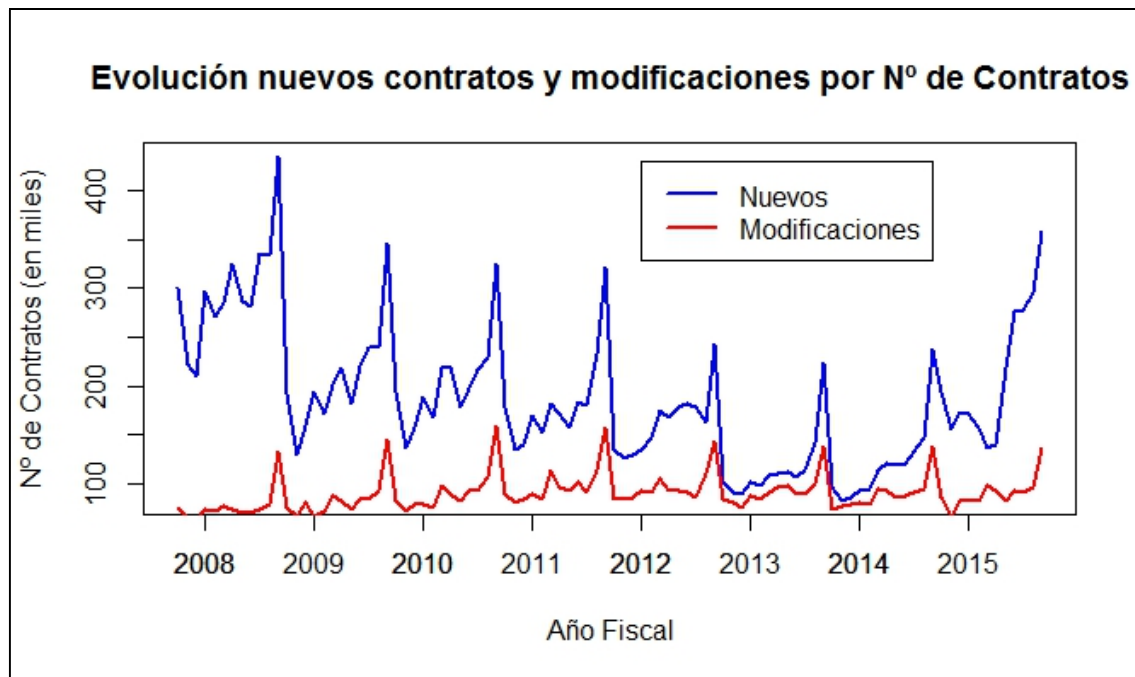


Figura 5.6 - Evolución del número de contratos vs. modificaciones para el período 2008/15

En la Figura 5.7, el comportamiento mostrado por la evolución del dinero presupuestado destinado a nuevas oportunidades registra una forma más regular que el que se dibujaba por número de nuevas oportunidades. Además, se puede apreciar como en muchos momentos el presupuesto destinado a modificaciones de contratos en curso es superior al que se destina a nuevas oportunidades.

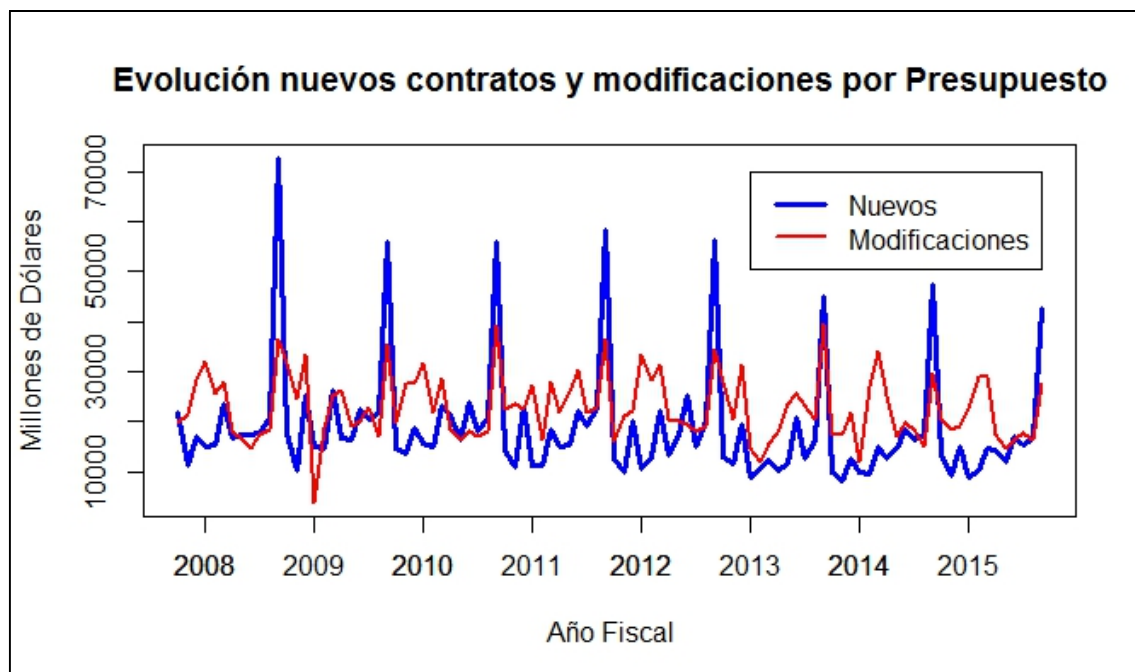


Figura 5.7 - Evolución del presupuesto de las nuevas oportunidades vs. modificaciones para el período 2008/15

5.7 Transformación de ficheros de nuevas oportunidades

En este punto se va a realizar la selección de las características definitivas para identificar los contratos y poder efectuar la predicción, transformando para ello los campos elegidos en valores numéricos discretos que permitan aplicar modelos matemáticos.

En Apache Spark, por ejemplo, se cuenta con multitud de clases en los paquetes de aprendizaje automático, *Spark MLlib* [51], que facilitan la extracción y el tratamiento de características empleadas en las tareas de aprendizaje. Clases como *QuantileDiscretizer* [48], que a partir de una columna con valores continuos genera una salida con valores categóricos agrupados. Los rangos para el agrupamiento se eligen tomando una muestra de los datos y se dividen en partes aproximadamente iguales, siendo los límites inferior y superior del total de los rangos $-\infty$ y $+\infty$, cubriendo así todos los valores reales, intentando determinar el mejor número de rangos para particionar la muestra de datos. La clase *Bucketizer* [49] mapea un conjunto de valores continuos a un conjunto de valores discretos predefinidos. Otras clases como *Normalizer* [50] o *StandardScaler* [52] ayudan a normalizar los datos o a escalarlos haciendo que el rango del dominio de valores donde toman valor los datos de entrada sea más reducido. De todas estas utilidades que se pueden encontrar en Spark a día de hoy, en Flink sólo está implementado el escalado, aunque se contempla en el *roadmap* del producto incluir próximamente otras [12].

La *discretización* de las características de los contratos a valores numéricos enteros se va a realizar utilizando una transformación de mapeo en la que a partir del literal de entrada de la característica correspondiente se le asigne un código numérico concreto.

En un mes y año concreto se puede saber si se ha dado un tipo de contrato con unas características determinadas o no, pudiendo tratar el problema como si fuera una clasificación entre si se da ese tipo de contrato o no. Tratarlo de esta manera engloba algunas dificultades, como describir todos los contratos que no se han dado para poder aplicar un aprendizaje supervisado. Además, para un mismo par mes-año se pueden dar varios contratos del mismo tipo y con las mismas características identificativas.

Se ha preferido, una vez seleccionadas algunas características de los contratos, agruparlos por ellas y ver cuántos se dan con esos parámetros en un mes y año concretos. De tal manera que se tiene una serie temporal desde el año fiscal 2008 para cada mes con el número de contratos con unas características determinadas que se han concedido. Para ello, una vez codificadas las características se va a agrupar por ellas.

Para esta primera fase de implementación se han seleccionado las siguientes características como definitivas, aunque en la fase de experimentación posterior se trabaje también con otras combinaciones distintas:

- Mes y Año (obtenidos del *signeddate*)
- Agencia (*agencyid*)
- Código NAICS (*principalnaicscode*)
- Rango de presupuesto (codificación del *dollarsobligated*)

Para la codificación en valores enteros de la agencia y del NAICS se va a partir de los ficheros resultantes de análisis realizado en R de las características en el apartado 5.2. En cada uno de estos ficheros aparece el listado total de los literales que pueden darse por cada característica y el número de contratos que tienen ese literal asignado. Es un listado ordenado de manera ascendente, así que el valor más repetido es el último del fichero. Lo que se va a hacer es asignar un entero según el orden que presenten en dicho fichero, manteniendo el 0 como reservado por si el campo no está informado (o está informado con un valor por defecto, como los dos puntos ':'). Por ejemplo, en el caso del NAICS, como el valor menos repetido es el “11: AGRICULTURE, FORESTRY, FISHING AND HUNTING” se le asignará el valor 1 y se irá incrementado en una unidad hasta llegar al más repetido, el “423450: MEDICAL, DENTAL, AND HOSPITAL EQUIPMENT AND SUPPLIES MERCHANT WHOLESALERS” al que se le asigna el valor 1.636. Los ficheros con el mapeo de la codificación aplicada a cada característica se generan en la clase Java encargada de la selección y transformación, *SeleccionConversion.java*.

Este tipo de transformación tiene sus riesgos y desventajas según el modelo de aprendizaje que después se aplique, pudiendo dar lugar a sesgos. Por ejemplo, se puede pensar que dos tipos de NAICS están más *cercanos* -en términos de similitud- si distan sólo una unidad en vez de 100, pero no es cierto: el ordinal que tiene asignado cada uno de ellos es por grado de repetición no por similitud en ningún otro término.

Para el caso de la discretización de los *dollarsobligated* se va a utilizar una función que aplique la codificación según los valores mostrados en la Tabla 5.2. Según el valor mostrado en el campo de *dollarsobligated*, se identifica el subrango de la tabla a la que pertenece esa cantidad y se genera como salida el número entero entre 0 y 15 que corresponda.

La clase que implementa la función de transformación de todas las características es la *SeleccionMapFunction.java*, que se emplea desde la clase *SeleccionConversion.java*. Esta clase se encarga primero de construir las tablas *hash* de las características a partir de los ficheros del análisis generado en R. Se leen esos ficheros y a cada valor se le asigna un entero; si la característica tenía el mismo código pero distinto texto (como en el caso de agencias y NAICS repetidos), no se registra en la tabla *hash* porque ya tiene entero asignado. En paralelo a la creación y poblado de esta tabla *hash* se van generando los ficheros de codificación donde queda representado que valor se le asigna a cada característica. Estas tablas se pasan a la clase de *SeleccionMapFunction* que será la encargada de aplicarlas sobre las características a las tuplas de los ficheros de contratos.

Una vez construidas las tablas *hash*, se van cargando cada uno de los ficheros de contratos de nuevas oportunidades en *datasets*. A todas las tuplas del *dataset* se les aplica la función de mapeo que transforma la agencia, el NAICS y el presupuesto, y que extrae el año y el mes: además se le añade un 1 como contador, para poder conocer el total en el agrupamiento. Se ha pasado de una tupla de 19 características de entrada a una tupla de 6: presupuesto, mes, año, NAICS, agencia y el contador.

Sobre esta nueva colección de tuplas se agrupa por mes, año, agencia, presupuesto y NAICS, sumando los contadores: así se obtiene el número de contratos concedidos en un mes concreto según la agencia responsable, el sector al que pertenece el proveedor y el montante económico estimado.

Se vuelca cada año fiscal en un fichero de salida con las características seleccionadas transformadas. En la Figura 5.8 puede verse una línea de ejemplo de uno de estos ficheros resultantes con la transformación realizada. Lo que indica esa línea del fichero es que en abril de 2014 la agencia “7008: U.S Coast Guard” concedió 11 contratos de entre tres y diez mil dólares a proveedores cuyo principal sector de dedicación es “332999: ALL OTHER MISCELLANEOUS FABRICATED METAL PRODUCT MANUFACTURING”.



Figura 5.8 - Esquema explicativo de una línea del fichero transformado

Para aquellas agrupaciones cuyo valor codificado del NAICS es 0 se eliminan. Carece de sentido mantenerlas teniendo en consideración que el objetivo es predecir cuántos contratos se pueden dar en un sector económico determinado que resulte de interés y si el código está a 0 es porque no estaba informado el valor en los registros de origen, haciendo inviable clasificar el tipo de servicio proporcionado por el proveedor adjudicatario del contrato.

5.8 Recuento

Desde el 2008 hasta el 2015, incluyendo ambos, son 8 años fiscales que deberían arrojar un total de 96 observaciones para tener la secuencia temporal completa. Si ocurre que para un par mes-año concreto no se ha adjudicado ningún contrato con unas características de agencia, presupuesto y NAICS para el que sí se ha adjudicado otros meses, la secuencia debería mostrar un 0 en el número de contratos concedidos en ese mes. En este punto del proceso eso no se está haciendo, pues se ha partido del fichero de contratos concedidos y sólo se reflejan los que han sido asignados, no hay ninguna referencia a aquellos que no han sido concedidos y que por tanto no existen.

Para poder completar las secuencias primero es necesario saber cuántas ternas de {presupuesto, NAICS, agencia} se han dado en el total de los años.

Eso se hace en la clase *Recuento.java*, que toma como origen los ficheros resultantes en el paso de selección y conversión, y donde se agrupan los registros por esos tres elementos.

Se van a generar dos ficheros de salida: uno con la colección completa de ternas {presupuesto, NAICS, agencia} que alguna vez se han dado, y otro fichero donde junto a cada terna se va a indicar el número de meses en los que han sido observados, para conocer del total de ternas y saber para cuántas están completas todas sus observaciones y para cuántas no.

El resultado de la ejecución de esta clase con esas características arroja un resultado de 202.046 ternas distintas. La distribución del número de observaciones puede observarse en el histograma de la Figura 5.9.

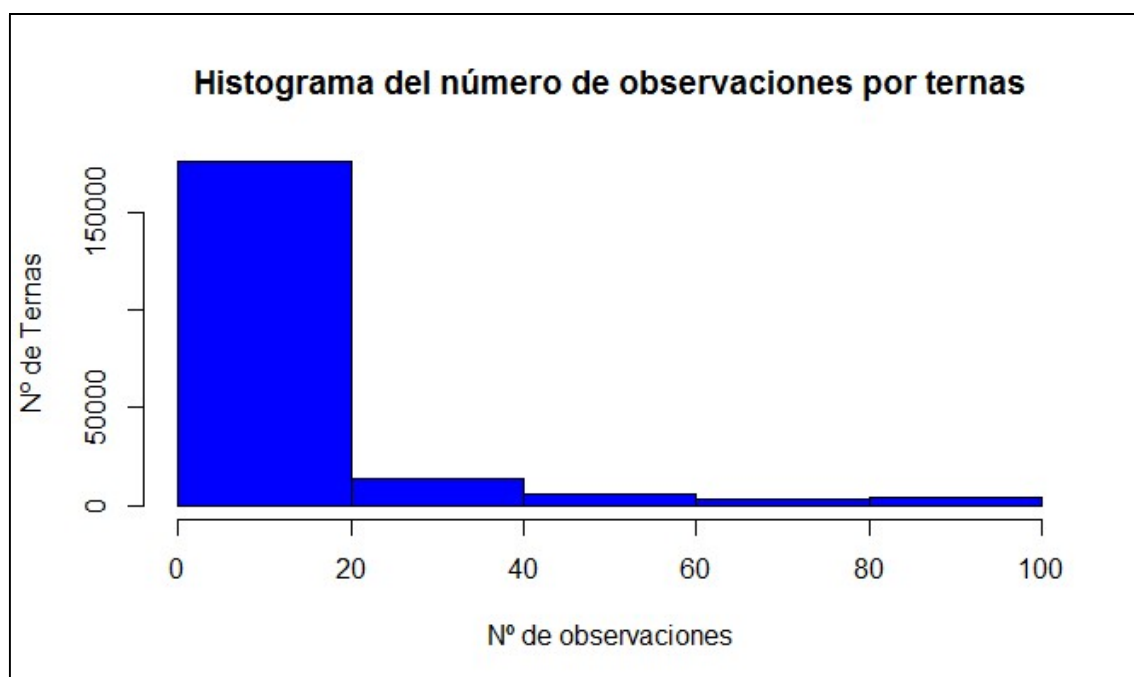


Figura 5.9 - Histograma del número de series con observaciones distintas de 0

De las más de doscientas mil totales, sólo 3.897 tienen 80 o más observaciones. 175.054 tienen menos de 20 observaciones. Sólo 1.130 tienen las 96 observaciones posibles (sin ningún mes en el que se hayan adjudicado 0 contratos) entre 2008 y 2015. Como resumen significativo: con más del 75% de observaciones rellenas sólo hay 4.961 ternas, que representan un 2,46% del total de ternas existentes.

5.9 Completación de la serie temporal

Como se indicaba antes, para un estudio correcto de la serie es necesario que todas las observaciones se hayan medido con la misma magnitud y utilizando la misma periodicidad para no alterar el sentido de la serie. Es en este punto del proceso es donde se va a realizar la completación de las series temporales, incluyendo un valor de 0 para aquellas observaciones de las ternas que no se hayan encontrado en la secuencia.

Este paso es especialmente importante para el estudio de la serie por métodos matemáticos, siendo imprescindible para garantizar un mínimo de integridad en los resultados que se obtengan sin verse desvirtuados por observaciones no representadas.

La clase *Completacion.java* toma como *dataset* de partida el fichero con las ternas encontradas generado en el paso anterior; además, se dispone de los ficheros de los contratos ya transformados y agrupados –los que contienen las observaciones mayores de cero–.

Para incluir las observaciones a 0 en cada uno de los ficheros de contratos se podría pensar en recorrer las ternas encontradas y, por cada año fiscal, con un bucle, comprobar si existe observación de la terna en cada mes y año, y de no existir incluir una nueva tupla en el *dataset* con los datos de la terna, el mes y año, y la cantidad de contratos a 0. Recorrer así todas las ternas (más de doscientas mil) y comprobar por cada una de ellas si existe en cada mes-año es un proceso largo y costoso. Apache Flink no ofrece ninguna ventaja adicional recorriendo secuencias en bucles, pero sí presenta un gran rendimiento tratando y agrupando grandes volúmenes de datos. Por ello, se ha descartado la idea del bucle y se ha partido de la creación de un fichero de plantilla válido para todos los años fiscales. Se carga el fichero de ternas encontradas en un *dataset* y se recorre una sola vez. Por cada terna y mes se incluye una línea en el fichero de plantilla: si el mes está entre octubre y diciembre se identifica el año como *\$year-1* y si está entre enero y septiembre como *\$year* –hay que recordar que el año fiscal de un año natural comienza en octubre del año anterior y finaliza en septiembre del año natural, de ahí las dos maneras diferenciadas de etiquetar el año en el fichero de plantilla–. Como número de contratos asociado a cada observación se fija un 0. De esta manera se genera un fichero de plantilla con 2.424.552 líneas: el resultado de multiplicar las 202.046 ternas por 12 observaciones anuales. En este punto, cada línea representa la observación de una terna para un mes concreto con 0 contratos concedidos.

Ahora, por cada año fiscal, se van a ir cargando uno a uno los ficheros de nuevas oportunidades ya transformados y el fichero de plantilla. Sobre este último, mediante una transformación implementada en la función de mapeo de la clase *MapeoAnyo.java*, se va a sustituir el *\$year* y *\$year-1* de las tuplas del fichero de plantilla por el valor que corresponda en ese año. Una vez transformado el *dataset* de las tuplas del fichero de plantilla se une al *dataset* de las tuplas del año fiscal que se esté evaluando en ese momento, y se aplica una transformación más (agregación): se agrupan las tuplas por presupuesto, mes, año, NAICS y agencia sumando el número de contratos. Las tuplas resultantes son aquellas que se buscaban y que se vuelcan al fichero de salida: tiene el número de contratos relleno para aquellos registros en los que había concesión de contratos y lo tienen a 0 (porque las ternas del fichero de plantilla no han encontrado con quién agruparse en la transformación) si no había concesiones de contratos.

Estos ficheros son los que se van a utilizar como materia prima de los procesos de aprendizaje y predicción posteriores.

5.10 Regresión lineal sobre múltiples predictores

Como primer modelo de aproximación se ha aplicado una regresión múltiple lineal para determinar el hiperplano que mejor se ajusta a la realidad determinada por las observaciones recogidas.

En esta primera fase de implementación se ha realizado la regresión utilizando como variables independientes las siguientes:

- Código numérico relativo al presupuesto
- Mes
- Año
- Código numérico que codifica al NAICS
- Código numérico que codifica la agencia

La variable dependiente de salida es el número de contratos concedidos.

El modelo de partida es completo en tanto que se están utilizando todas las variables que quedaron definidas en los ficheros generados para la construcción del modelo. En este caso no es posible añadir directamente más variables como predictores directamente –siguiendo la estrategia de *forwarding* o inclusión de predictores- porque sería necesario volver a procesar los ficheros y agrupar el número de contratos según las nuevas características de los contratos elegidas. Lo que si es factible hacer es seguir una estrategia de *back elimination* (eliminación de predictores), donde, si se detecta algún predictor redundante o que no aporta mejoría al modelo, se pueda prescindir de él. En el apartado de experimentación se incluyen otros modelos de regresión construidos utilizando conjuntos distintos de características de los contratos como variables independientes.

Para la construcción del modelo en este punto, desde R Studio se cargan los ficheros con las series completas para todos los años fiscales entre 2008 y 2014. El año 2015 se va a utilizar para la evaluación del modelo generado estudiando así la calidad con la que se puede aproximar 2015 a partir de los datos del período 2008-2014.

Para la construcción del modelo de regresión se han seguido estos pasos:

- **Paso 1.** Carga de los ficheros de datos y construcción del *dataframe* con todas las observaciones. Al *dataframe* se le ha llamado ‘*contratos*’ y el objeto *vars* contiene el nombre de las columnas de esta forma:

```
vars <-  
c('presupuesto','mes','año','naics','agencia','cantidad_contratos')
```

- **Paso 2.** Cálculo de los coeficientes de correlación

Se puede realizar de esta forma:

```
cor(contratos[,vars], use="pairwise.complete.obs")
```


O de esta utilizando el paquete *psych*:

```
library(psych)
corr.test(contratos[,vars], use="pairwise.complete.obs")
```

Con la función *corr.test*, los coeficientes de correlación obtenidos y redondeados a 2 decimales se muestran en la Figura 5.10. Hay 16.971.864 de registros en la muestra de datos empleada para la elaboración del modelo. No se aprecia ninguna correlación significativa de variables no llegando ninguna al 0,2 en valor absoluto.

Correlation matrix							
	presupuesto	mes	año	naics	agencia	cantidad_contratos	
presupuesto	1.00	0.00	0.00	0.08	0.04	-0.01	
mes	0.00	1.00	-0.16	0.00	0.00	0.00	
año	0.00	-0.16	1.00	0.00	0.00	0.00	
naics	0.08	0.00	0.00	1.00	-0.15	0.01	
agencia	0.04	0.00	0.00	-0.15	1.00	0.02	
cantidad_contratos	-0.01	0.00	0.00	0.01	0.02	1.00	
Sample Size							
[1]	16971864						

Figura 5.10 - Valor de los coeficientes de correlación entre las variables independientes a emplear en el modelo de regresión

- **Paso 3.** Se aplica la función *lm* para el cálculo del modelo y se visualizan los coeficientes. Para ello se realiza:

```
Modelo_lm <- lm(cantidad_contratos ~
presupuesto+mes+año+naics+agencia, data = contratos)
summary(Modelo_lm)
```

Cuando se utilizan todo el resto de variables como predictores puede indicarse así la sentencia para generar el modelo:

```
Modelo_lm <- lm(cantidad_contratos ~ ., data = contratos)
```

Los resultados detallados del modelo obtenido se encuentran en la Figura 5.11. Conociendo el valor de los coeficientes que ofrece el modelo:

- Interceptor: 152.328929471
- Dinero: -0.151590031
- Mes: -0.005645934
- Año: -0.078491461
- Naics: 0.001860133
- Agencia: 0.019849267

Se puede formular el número de contratos (redondeando decimales) según la fórmula:

$$N^{\circ} \text{ Contratos} = 152.329 - 0.152 * \text{Dinero} - 0.006 * \text{Mes} - 0.078 * \text{Año} + 0.002 * \text{Naics} + 0.02 * \text{Agencia}$$

```

Residuals:
    Min       1Q   Median       3Q      Max
   -3.2    -1.5    -0.9    -0.1   17836.1

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.523e+02  9.681e+00  15.734  <2e-16 ***
presupuesto -1.516e-01  2.970e-03 -51.039  <2e-16 ***
mes          -5.646e-03  2.853e-03  -1.979   0.0478 *
año          -7.849e-02  4.813e-03 -16.308  <2e-16 ***
naics         1.860e-03  2.710e-05   68.651  <2e-16 ***
agencia       1.985e-02  2.415e-04   82.198  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 40.06 on 16971858 degrees of freedom
Multiple R-squared:  0.000708, Adjusted R-squared:  0.0007077
F-statistic: 2405 on 5 and 16971858 DF, p-value: < 2.2e-16

```

Figura 5.11 - Parámetros del modelo de regresión múltiple lineal generado

El coeficiente de determinación ajustado ofrecido por el modelo, R^2_{adj} , es 0,00071, lo que indica que sólo un 0,07% de la varianza de los datos se explica con el modelo obtenido.

- **Paso 4.** Se calcula el Factor de Inflación de la Varianza (VIF), con la función `vif(Modelo_lm)`, para ver si alguna de las variables supera ese umbral de 5 que las señala como redundantes. Los valores obtenidos se representan en la Figura 5.12. Como ninguna variable lo supera se mantienen todas como predictores.

presupuesto	mes	año	naics	agencia
1.010129	1.026023	1.026023	1.030088	1.024954

Figura 5.12 - Resultados para la función `vif(Modelo_lm)`

El siguiente paso a seguir es la representación gráfica que ayude a demostrar si el modelo obtenido cumple las premisas que debe seguir uno adecuado de regresión lineal. Para ello, el primer gráfico se muestra en la Figura 5.13 presentando la relación entre los valores ajustados y los residuos.

Si el modelo se ajustara de manera homogénea a los datos de entrada este gráfico debería ser uniforme. La acumulación de valores en un extremo del gráfico denota que existe heterocedasticidad y que el modelo no ajusta por igual todos los valores.

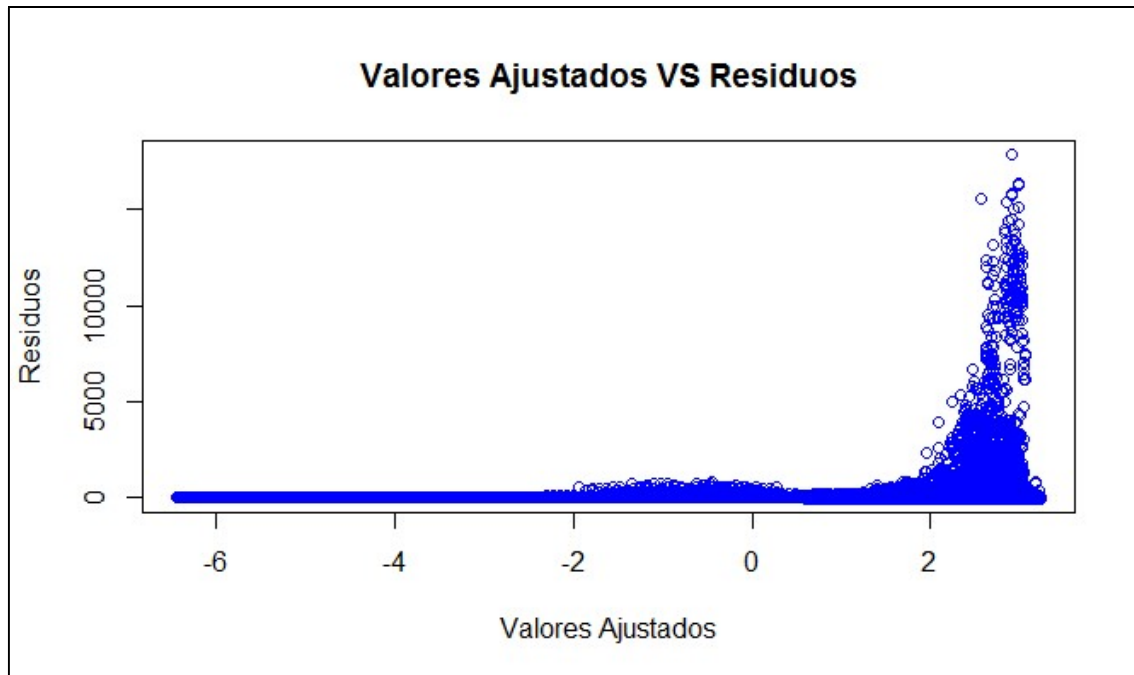


Figura 5.13 - Gráfico que contrasta los valores ajustados frente a los residuos

Para generar los dos siguientes gráficos también se ha utilizado la función *plot*, empleada en la generación del gráfico anterior. Al incluir el parámetro *which* para seleccionar el tipo de gráfico asociado al modelo de *lm* (*linear model*) se hace por parte de R un uso interno de la función *plot.lm* que tiene codificados estos 5 gráficos: "*Residuals vs Fitted*", "*Normal Q-Q*", "*Scale-Location*", "*Cook's distance*", "*Residuals vs Leverage*". El pequeño inconveniente que presenta su uso es que todos los textos que se muestran en los gráficos están *harcodados* en la función de generación de los gráficos y no son parametrizables para traducirlos al español y poder mostrarlos adecuadamente junto al texto de esta memoria de trabajo. Para solventarlo, lo que se ha hecho es copiar la función *plot.lm* y editarla para modificar el título de los gráficos y las etiquetas asociadas a los ejes vertical y horizontal. Para copiar la función y editarla hay que utilizar las siguientes instrucciones:

```
plotlm <- stats::plot.lm
fix("plotlm")
```

Otro de los gráficos generados para comprobar la idoneidad del modelo de regresión es el gráfico Normal Q-Q ("Q" viene de cuantil), donde se muestra la relación entre los cuantiles de los residuos y los cuantiles teóricos. El gráfico elaborado lo constituye la Figura 5.14. Para que el modelo fuera idóneo los residuos deberían mostrarse ajustados a la línea roja que aparece en el gráfico y no mostrar ningún patrón de agrupamiento.

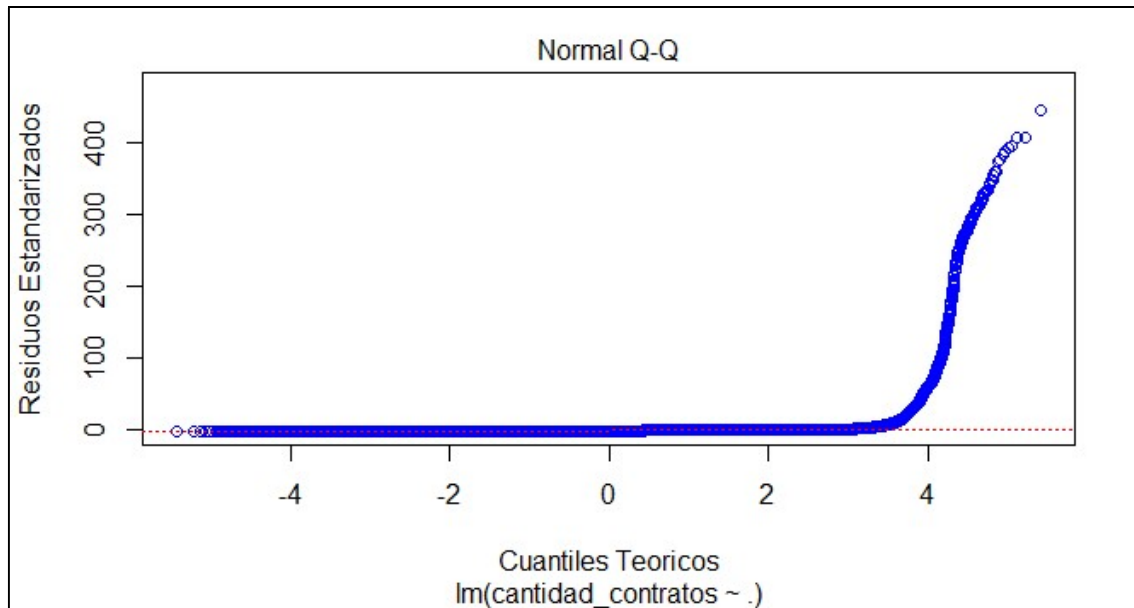


Figura 5.14 - Gráfico Normal Q-Q

El tercer y último gráfico generado es el de distancias de Cook mostrado en la Figura 5.15. Se han identificado las seis observaciones que más peso tienen en la generación del modelo.

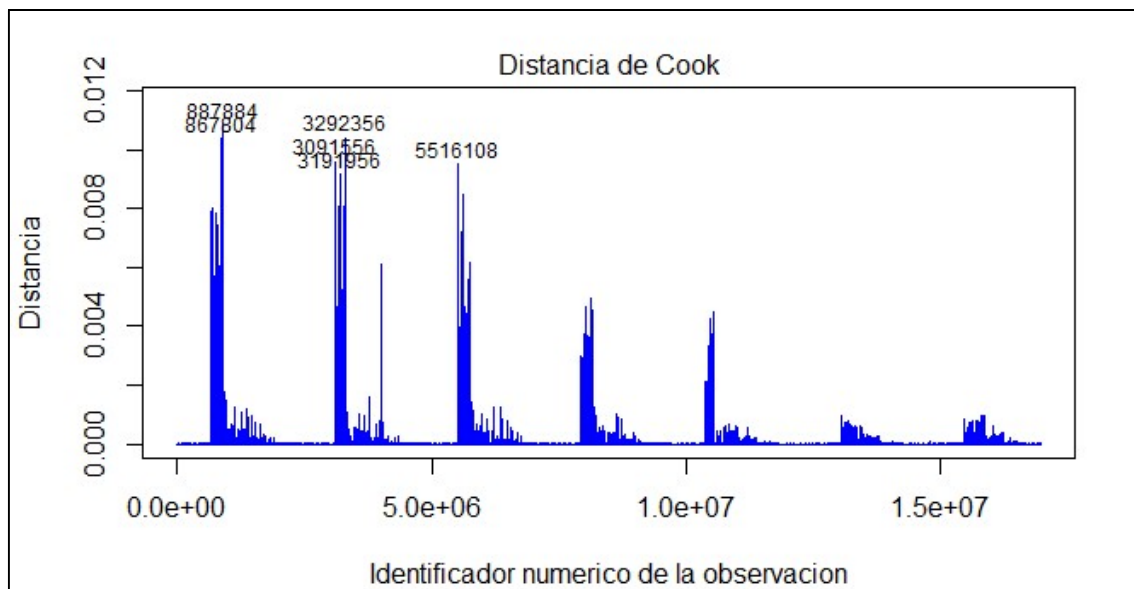


Figura 5.15 - Gráfico de Distancias de Cook

Los valores señalados por el gráfico de distancias de Cook dentro del *dataframe* de contratos son:

```
>contratos[887884,]
      presupuesto mes  año naics agencia cantidad_contratos
887884           3   9 2008  1631      289             16385
> contratos[867804,]
      presupuesto mes  año naics agencia cantidad_contratos
867804           3   8 2008  1631      289             16300
> contratos[3292356,]
      presupuesto mes  año naics agencia cantidad_contratos
```

```

3292356      3      8 2009  1631      289      17839
> contratos[3091556,]
      presupuesto mes año naics agencia cantidad_contratos
3091556      3      1 2009  1631      289      13722
> contratos[3191956,]
      presupuesto mes año naics agencia cantidad_contratos
3191956      3      3 2009  1631      289      15041
> contratos[5516108,]
      presupuesto mes año naics agencia cantidad_contratos
5516108      3      1 2010  1631      289      14430

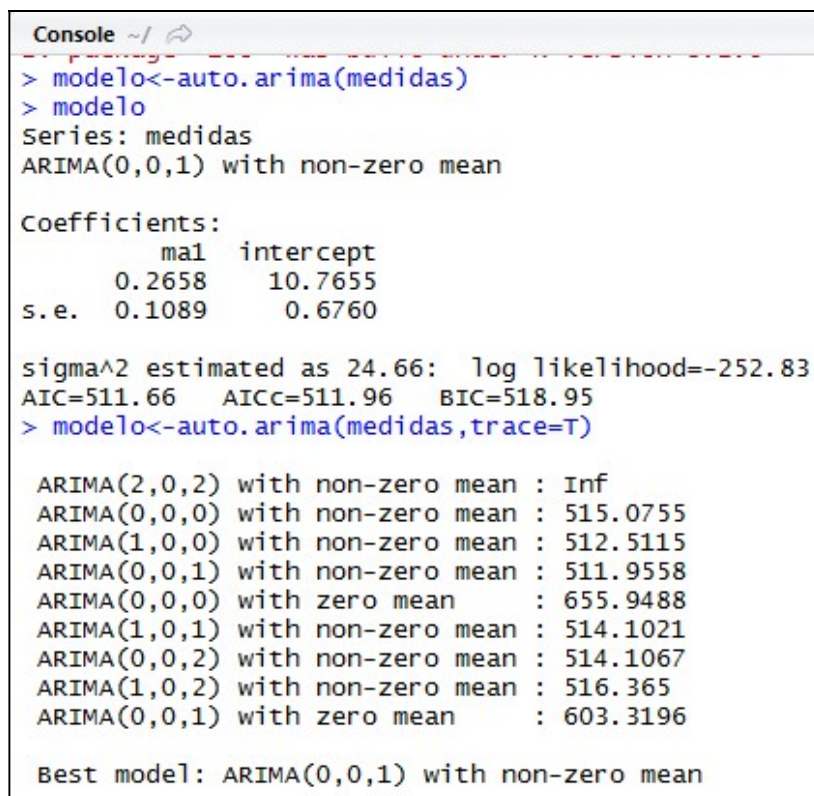
```

Las observaciones que contienen un número de contratos más alto, como estas seis mostradas que están por encima de los trece mil mensuales han tenido mucho más peso en la generación que el resto.

5.11 Función `auto.arima` sobre series concretas

La función *auto.arima* se aplica de manera individual a cada serie concreta representada por las observaciones registradas para una terna particular de {presupuesto, NAICS, agencia}. A diferencia del modelo de regresión que se ha construido utilizando todos los datos de todas las series, en este caso es necesario conseguir la serie ordenada de las 96 observaciones para cada una de las más de doscientas mil ternas que se habían encontrado.

Como ejemplo de implementación se puede ver en la Figura 5.16 la salida en la consola de R Studio de la aplicación de la función *auto.arima* a una de las series de contratos.



```

Console ~/
> modelo<-auto.arima(medidas)
> modelo
Series: medidas
ARIMA(0,0,1) with non-zero mean

Coefficients:
      ma1  intercept
      0.2658    10.7655
s.e.    0.1089     0.6760

sigma^2 estimated as 24.66:  log likelihood=-252.83
AIC=511.66  AICc=511.96  BIC=518.95
> modelo<-auto.arima(medidas,trace=T)

ARIMA(2,0,2) with non-zero mean : Inf
ARIMA(0,0,0) with non-zero mean : 515.0755
ARIMA(1,0,0) with non-zero mean : 512.5115
ARIMA(0,0,1) with non-zero mean : 511.9558
ARIMA(0,0,0) with zero mean : 655.9488
ARIMA(1,0,1) with non-zero mean : 514.1021
ARIMA(0,0,2) with non-zero mean : 514.1067
ARIMA(1,0,2) with non-zero mean : 516.365
ARIMA(0,0,1) with zero mean : 603.3196

Best model: ARIMA(0,0,1) with non-zero mean

```

Figura 5.16 - Ejemplo aplicación función *auto.arima*

Los valores detallados del modelo ARIMA(0,0,1) obtenido que mejor ajustaba los datos de la serie se muestran en la Figura 5.17. Los tres primeros son los obtenidos por los criterios que determinan que ese modelo presenta el mejor ajuste (por defecto se aplican por la función los tres criterios –AIC, AICc y BIC- y el valor resultante es el del menor de los tres encontrados). El cuarto hace referencia a los coeficientes del modelo ARIMA(p,d,q): el modelo que mejor ajusta es puramente de media móvil, MA, de orden 1 (su valor de q vale 1).

```
> modelo$aic
[1] 511.6558
> modelo$aicc
[1] 511.9558
> modelo$bic
[1] 518.9482
> modelo$coef
      ma1 intercept
0.2658334 10.7655415
```

Figura 5.17 - Valores del modelo ARIMA de los criterios y de los coeficientes del resultado

Aplicando la función de predicción sobre el modelo obtenido para las siguientes 12 observaciones –que serían las correspondientes al año fiscal 2015- se obtienen los valores de la Figura 5.18. Esos valores se han enfrentado a los valores reales ofrecidos por la serie en 2015 y se pueden consultar en la Tabla 5.4.

```
> resultado$mean
Time series:
start = 85
End = 96
Frequency = 1
[1] 12.95610 10.76554 10.76554 10.76554 10.76554 10.76554
[7] 10.76554 10.76554 10.76554 10.76554 10.76554 10.76554
```

Figura 5.18 - Valores de la predicción para 12 observaciones con ARIMA(0,0,1) obtenido previamente a partir de las anteriores observaciones de la serie

Presupuesto	Mes	Año	Naics	Agencia	N.º Contratos Real	ARIMA (0,0,1)
9	10	2014	202	287	7	12,96
9	11	2014	202	287	12	10,77
9	12	2014	202	287	12	10,77
9	1	2015	202	287	14	10,77
9	2	2015	202	287	10	10,77
9	3	2015	202	287	12	10,77
9	4	2015	202	287	15	10,77
9	5	2015	202	287	13	10,77
9	6	2015	202	287	15	10,77
9	7	2015	202	287	14	10,77
9	8	2015	202	287	14	10,77
9	9	2015	202	287	25	10,77

Tabla 5.4 - Serie de datos real y aproximada por ARIMA para el año 2015 de la terna {presupuesto:9, NAICS:202, agencia:287}

La aplicación de esta misma función desde una clase Java que haga uso de la interfaz JRI con R se hace con las líneas de código mostradas en la Figura 5.19. Previamente se ha guardado en un fichero de texto los datos de la serie con las observaciones del número de contratos. Se incluye la implementación de la carga del fichero, aplicación de *auto.arima* y predicción de las 12 observaciones siguientes.

```
REXP x;
re.eval("medidas <- read.csv(\"C:/PFC/workspace/c9_202_287.txt\", \"+
    \"stringsAsFactors=FALSE)\");
re.eval("names(medidas) <- c('orden', 'serie')");
re.eval("medidas <- medidas$serie");
re.eval("library(zoo)");
re.eval("library(timeDate)");
re.eval("library(forecast)");
x = re.eval("medidas");
x = re.eval("modelo<-auto.arima(medidas,trace=T)");
x = re.eval("modelo");
re.eval("resultado<-forecast(modelo,h=12)$mean");
```

Figura 5.19 - Extracto de ejemplo de implementación en Java del uso de *auto.arima* en R

5.12 Redes LSTM

El modelo implementado de red neuronal artificial consta de cuatro capas: las tres primeras capas son neuronas de tipo LSTM y la última capa, la de salida, es de tipo prealimentada. En la primera capa, la de entrada, se tiene una única neurona con salida a las 35 neuronas de la capa siguiente. Esta segunda capa, a su vez, comunica con las 70 neuronas de una segunda capa LSTM oculta. A la capa de 35 y 70 neuronas se le aplica un *dropout* de 0.15 en cada una de ellas. Por último, todas las salidas de la tercera capa se conectan con las doce neuronas de la capa de salida (todas conectadas, capa tipo Dense), para la que se fija una función de activación lineal.

La pérdida del modelo se evalúa según el error cuadrático medio (MSE, cuya fórmula se incluye más adelante) calculado y el optimizador elegido es el *RMSprop*, es recomendado para redes recurrentes.

El número de épocas utilizado es 4 y el tamaño de cada lote empleado antes de ajustar pesos es de 512 elementos de entrada, fijando una proporción de datos de validación del 20%.

La fórmula para el ECM (en inglés, *MSE*, *Mean Squared Error*) es la siguiente:

$$ECM = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

El esquema de la red se ha representado en la Figura 5.20.

□ Parámetros red neuronal LSTM en Keras

Para implementar la red LSTM anteriormente descrita en Keras utilizando TensorFlow como *backend* es necesario definir los siguientes parámetros para la implementación y entrenamiento de la red:

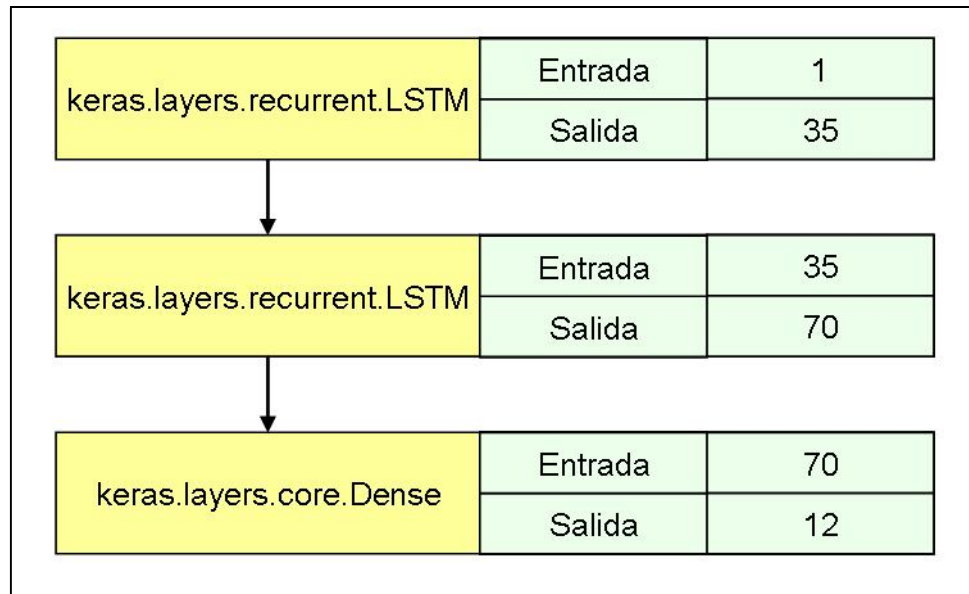


Figura 5.20 - Esquema red de neuronas implementado en Keras/TensorFlow

- Número de capas y la cantidad de celdas LSTM que va a contener cada capa. El modelo se define como *Sequential()* que indica que se van a ir añadiendo al modelo cada una de las capas, una tras otra.
- Cada capa a incluir en el modelo será de tipo *LSTM()*. Para cada una de ellas se indica:
 - *input_dim*: el número de entradas que recibe esta capa.
 - *output_dim*: el número de salidas que genera esta capa.
 - *return_sequences*: *True* o *False*. Indica si la información de salida de una capa debe transferirse sólo a la siguiente capa (*False*) o si además de transferirse a ésta debe utilizarse en ella misma en la siguiente iteración (*True*).
 - *weights*: si se quieren inicializar los pesos de la red.
 - *input_length*: es obligatorio si se utiliza TensorFlow junto con Keras. Sirve para indicar el número de campos –columnas- que contiene la serie de datos de entrada (84, en este caso, que son las observaciones de entrenamiento entre 2008 y 2014).
- A una capa se le puede añadir un elemento de *dropout* (retiro, desactivación). *Dropout* [58] es una técnica para evitar el sobreajustamiento (*overfitting*) de la red a los datos de entrada y que consisten en la desactivación de una parte aleatoria de los elementos de la red durante cada fase de actualización. El *dropout* se fija como un valor entre 0 y 1, indicando la proporción de elementos a retirar. En este caso se ha empleado 0,15.
- Para finalizar la red se ha incluido una última capa, pero esta vez no es de tipo *LSTM*, es una capa prealimentada normal donde aplicar la regresión a los valores obtenidos en la secuencia anterior. Para ello se utiliza una de tipo *Dense*, que modela una capa donde todas las neuronas están conectadas con todas (conexión

completa). A esta capa se le añade una función de activación lineal para facilitar la regresión (existen otras activaciones disponibles como *softmax*, *softplus*, *sigmoid*, etc.).

A la hora de compilar el modelo es necesario indicar varias cosas:

- El modo de calcular la pérdida (*loss*) entre el valor esperado y el obtenido. Es la función objetivo y para una regresión el típico es el *MSE* (*Mean Squared Error*, Error Cuadrático Medio), pero pueden indicarse otros como: *MAE* (*Mean Absolute Error*), *MSLE* (*Mean Squared Logarithmic Error*), *poisson*; *categorical_crossentropy*, etc.).
- El método de optimización, por el cuál se ajustan los pesos buscando la minimización de la función objetivo. Se puede indicar *SGD* para el uso del gradiente estocástico descendente, *Adam*, *Adagrad*, o *RMSprop* que es el recomendado en redes recurrentes (*RMSprop* se define a partir de tres parámetros, que por defecto en Keras son el *Learning Rate* (Ratio de Aprendizaje): *lr*=0.001, *rho*=0.9, y un Factor de Incertidumbre: *epsilon*=1e-08).

En el momento de entrenar la red se definen los siguientes parámetros:

- ***x, y***: conjunto de datos de test y conjunto de salidas esperadas, respectivamente.
- ***nb_epoch***: el número de veces (épocas) que se ha de emplear toda la colección de datos de entrada para entrenar la red.
- ***verbose***: nivel de trazabilidad (0, 1, o 2). Siendo 0: sin trazas, 1: todas las trazas, y 2: una línea de traza por cada época.
- ***batch_size***: es el número de elementos de la colección de datos a utilizar en cada fase de entrenamiento antes de ajustar de nuevo los pesos. A mayor *batch_size*, mayor necesidad de memoria. Para cumplir una época se ha de haber empleado toda la colección de datos, utilizando en cada iteración un *batch_size*.
- ***validation_split***: proporción de los datos de entrada que se van a utilizar en la validación del modelo en la fase de entrenamiento. Se aparta esta porción de los datos de entrada, no se utilizan en el entrenamiento y después se utilizan para evaluar la pérdida y las métricas del modelo al final de cada época.
- ***shuffle***: determina si se quieren mezclar los datos de entrada o no, para no entrenar la red siempre en el mismo orden de aparición de estos.

En la Figura 5.21 se incluye el código de la implementación final de la red empleada en esta fase, donde se pueden ver los valores que han tomado los parámetros descritos anteriormente.

```
def construir_modelo():

    model = Sequential()
    layers = [1, 35, 70, 12]

    model.add(LSTM(
        input_dim=layers[0],
        output_dim=layers[1],
        return_sequences=True, input_length=84, input_shape=(84, 12)))
    model.add(Dropout(0.15))

    model.add(LSTM(
        layers[2],
        return_sequences=False))
    model.add(Dropout(0.15))

    model.add(Dense(
        output_dim=layers[3]))
    model.add(Activation("linear"))

    start = time.time()
    model.compile(loss="mse", optimizer="rmsprop")
    print("Tiempo de compilación : ")
    print(time.time() - start)
    return model
```

Figura 5.21 - Implementación en Python empleando Keras

La principal ventaja que ofrecía Keras/TensorFlow era la rapidez en el diseño del modelo y en su experimentación. Antes de llegar a esta implementación definitiva se han testado otras variantes con resultados dispares eligiendo el modelo presentado anteriormente por ser el que menor error presentaba en la evaluación de la pérdida de los datos validación.

Las trazas obtenidas durante las pruebas para cada una de las alternativas se muestran a continuación. En cada traza se muestra el número de elementos del conjunto de entrenamiento no destinados a validación, el tiempo empleado en la época, y el valor de pérdida para el conjunto de entrenamiento y de pérdida para el conjunto de validación.

Alternativa A) Arquitectura 1-40-80-12 neuronas, 3 épocas y un 10% datos para validación

```
Entrenamiento del modelo
Train on 172748 samples, validate on 9093 samples
Epoch 1/3
172748/172748 [====] - 1123s - loss: 9215.3132 - val_loss: 1072.1098
Epoch 2/3
172748/172748 [====] - 1107s - loss: 9205.4263 - val_loss: 1066.3948
Epoch 3/3
172748/172748 [====] - 1101s - loss: 9198.2509 - val_loss: 1061.5939
```

Tiempo total empleado: 55,5 minutos. Error validación medio: 1.066,70.

Alternativa B) Arquitectura 1-45-90-12 neuronas, 4 épocas y un 10% datos para validación

```
Entrenamiento del modelo
Train on 172748 samples, validate on 9093 samples
Epoch 1/4
172748/172748 [====] - 1402s - loss: 8964.7734 - val_loss: 2431.4249
Epoch 2/4
172748/172748 [====] - 1368s - loss: 8953.2524 - val_loss: 2418.1920
Epoch 3/4
172748/172748 [====] - 1358s - loss: 8945.1568 - val_loss: 2409.2324
Epoch 4/4
172748/172748 [====] - 1314s - loss: 8938.9203 - val_loss: 2399.8968
```

Tiempo total empleado: 90,7 minutos. Error validación medio: 2.414,69.

Alternativa C) Arquitectura 1-35-70-12 neuronas, 4 épocas y un 20% datos para validación

```
Entrenando el modelo
Train on 153555 samples, validate on 8082 samples
Epoch 1/4
153555/153555 [====] - 1002s - loss: 9874.4136 - val_loss: 58.7036
Epoch 2/4
153555/153555 [====] - 972s - loss: 9865.7964 - val_loss: 56.6883
Epoch 3/4
153555/153555 [====] - 1015s - loss: 9862.7066 - val_loss: 57.2934
Epoch 4/4
153555/153555 [====] - 974s - loss: 9860.7208 - val_loss: 53.9954
```

Tiempo total empleado: 60 minutos. Error validación medio: 56,67.

Alternativa D) Arquitectura 1-45-80-12 neuronas, 4 épocas y un 20% datos para validación

```
Entrenando el modelo
Train on 153555 samples, validate on 8082 samples
Epoch 1/4
153555/153555 [====] - 1134s - loss: 9712.9171 - val_loss: 167.0851
Epoch 2/4
153555/153555 [====] - 1201s - loss: 9707.8097 - val_loss: 164.2334
Epoch 3/4
153555/153555 [====] - 1178s - loss: 9703.0299 - val_loss: 161.0425
Epoch 4/4
153555/153555 [====] - 1114s - loss: 9698.3858 - val_loss: 157.7040
```

Tiempo total empleado: 77,12 minutos. Error validación medio: 162,52.

6. Experimentación

6.1 Rendimiento en el tratamiento de datos

En este punto se muestran los tiempos de procesamiento utilizados en cada una de las fases del tratamiento de los ficheros de contratos. Para ello hay que tener en cuenta que el dispositivo donde se han realizado todas las ejecuciones ha sido el portátil *Toshiba Satellite Pro C70-B-34T*, *Intel Core i5-5200U 2.20GHz*, que cuenta con 16 GB de memoria RAM y 4 procesadores lógicos.

Los tiempos empleados se han medido siempre en las mismas condiciones de carga de trabajo del equipo -sin ningún otro proceso reseñable en curso- y se incluye una comparativa de cuánto se ha consumido utilizando o no paralelismo en la ejecución.

Primero, en la Tabla 6.1, se presentan los datos individuales del tiempo de procesamiento de un solo fichero, el del año 2014, para la fase de Preselección y para la fase de División. Se utiliza un paralelismo de escritura en fichero de índice 1 en ambos escenarios, para garantizar así que el fichero de salida es uno solo y no varios. Como se observa en los resultados presentados, el tiempo empleado en las transformaciones utilizando un índice de paralelización mayor es mejor, consumiendo menos tiempo.

Fase	Paralelismo Transformaciones	Tiempo empleado
Preselección	4	93 segundos
	1	144 segundos
División	4	43 segundos
	1	49 segundos

Tabla 6.1 - Rendimientos para un sólo fichero

Los datos para el tratamiento de la totalidad de los ficheros se recogen en la tabla Tabla 6.2.

Fase	Paralelismo Transformaciones	Tiempo empleado
Preselección	4	92 minutos 31 segundos
	1	67 minutos 30 segundos
División	4	60 minutos 59 segundos
	1	56 minutos 30 segundos
Selección y Conversión	4	2 minutos 48 segundos
	1	2 minutos 49 segundos
Recuento	4	21 segundos
	1	30 segundos
Completación	4	7 minutos 23 segundos
	1	6 minutos 53 segundos

Tabla 6.2 - Rendimiento para la totalidad de los ficheros

En el cálculo del rendimiento de la totalidad de los ficheros también se ha empleado un índice de paralelización en escritura de 1. Puede observarse que al tratarse todos los ficheros con un índice de paralelización mayor de uno y utilizando un bucle *for* para su recorrido el tiempo empleado es mayor que con un índice de paralelización menor, especialmente en las fases de Preselección y División. En estas fases se están tratando los 16 ficheros, del año 2000 al 2015 incluidos. Esto ocurre porque al utilizar paralelización superior con varios ficheros se aplican las transformaciones en paralelo pero de múltiples ficheros a la vez que no se pueden escribir en disco hasta no estar todas las particiones de un mismo fichero procesadas, lo que puede provocar esperas que aumentan el tiempo total de ejecución. Para mejorar el rendimiento y aprovechar plenamente la capacidad de paralelización lo mejor sería lanzar trabajos por separado con un índice de paralelización de 4.

Los tiempos en las fases posteriores –Selección y Conversión, Recuento y Completación– son bastante similares, independientemente del índice de paralelización empleado. En estas fases ya sólo se emplean 8 ficheros, de los años 2008 y 2015, por lo que no se aprecia la penalización que se sufría en las fases anteriores con un índice de paralelización 4.

6.2 Experimentación agrupando por presupuesto y NAICS

Buscando comprobar si ocurre una mejoría significativa de la calidad del modelo de regresión múltiple utilizando colecciones de predictores distintas, se han realizado dos experimentos más que requieren nuevas clases Java para seleccionar y agrupar por las nuevas características escogidas. La arquitectura de estas nuevas clases puede observarse en el diagrama de clases de la Figura 6.1, conteniendo los sufijos *Experimentacion2* y *Experimentacion3* aquellas que son variaciones de las clases presentadas en la fase de implementación y que sólo varían algunas líneas de sus métodos. El resto de clases se explican en el apartado donde se trata la comparativa de los distintos modelos.

El modelo de regresión lineal múltiple utilizando los datos agrupados por la terna {presupuesto, NAICS, agencia} no explicaba un porcentaje aceptable de la varianza que reflejaban los datos. En esta nueva experimentación se va a comprobar cómo evoluciona la calidad del modelo eliminando el predictor de *agencia* y conservando los otros dos. Para ello es necesario implementar una nueva versión de las clases *SeleccionConversion.java*, *Recuento.java* y *Completacion.java* que reflejen la agrupación de contratos con los nuevos criterios; como punto de partida se utilizan los mismos ficheros de nuevas oportunidades que ya se habían generado antes.

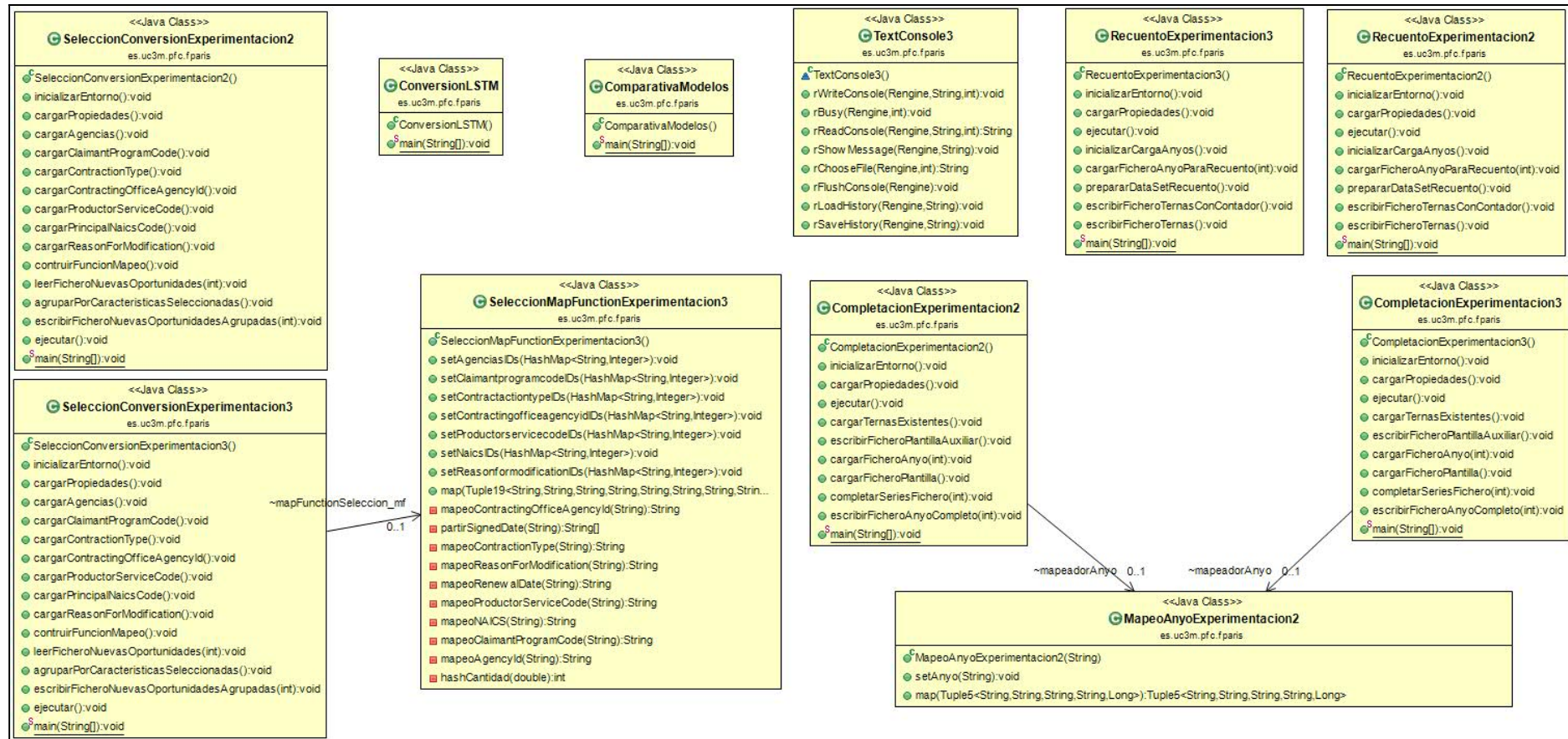


Figura 6.1 - Diagrama de Clases empleadas en la Experimentación

Los resultados de aplicar el modelo de regresión a las agrupaciones de contratos por NAICS y presupuesto se incluyen en la Figura 6.2.

Residuals:					
Min	1Q	Median	3Q	Max	
-56.3	-17.5	-7.6	3.8	17891.7	
Coefficients:					
	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	2.055e+03	1.313e+02	15.654	<2e-16	***
dinero	-2.837e+00	3.406e-02	-83.277	<2e-16	***
mes	-7.375e-02	3.869e-02	-1.906	0.0566	.
year	-1.025e+00	6.527e-02	-15.709	<2e-16	***
naics	3.938e-02	3.522e-04	111.826	<2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1					
Residual standard error: 150.3 on 1299307 degrees of freedom					
Multiple R-squared: 0.01321, Adjusted R-squared: 0.01321					
F-statistic: 4349 on 4 and 1299307 DF, p-value: < 2.2e-16					

Figura 6.2 - Regresión múltiple empleando sólo presupuesto y NAICS

El valor del coeficiente de determinación ajustado sube hasta el 0,01321, por lo que poco más del 1,3% de la varianza de los datos es explicable con este modelo. Se ha duplicado desde el 0,7 anterior pero aún resulta insuficiente teniendo en cuenta que se está prescindiendo de la información de la *agencia*, que podría ser de utilidad a un contratista futuro. Los gráficos realizados para comprobar la distribución del modelo siguen presentando los mismos problemas que los explicados en el apartado de implementación: heterocedasticidad y gran influencia de los valores numéricos muy altos en el modelo.

6.3 Experimentación agrupando por presupuesto y reduciendo el NAICS

Este tercer intento mantiene como características el presupuesto y el NAICS, que son las mínimas que pueden describir un contrato desde el punto de vista necesario para cubrir el objetivo perseguido en este proyecto. Hasta ahora para el NAICS se empleaba el código numérico completo que identifica unívocamente un sector, pero es posible agrupar varios sectores en categorías más genéricas a partir de las dos primeras cifras de este código.

Se distinguen 20 categorías económicas principales para el NAICS: si se reducen los 1.637 códigos de la primera experimentación a estas dos decenas, el número de contratos por mes estará menos disperso y serán menos las series que tengan observaciones a 0.

En la Tabla 6.3 se incluye la codificación de sectores de actividad por código NAICS reducido.

Código del sector	Descripción de la actividad
11	Agricultura, silvicultura, caza y pesca
21	Minería, canteras, extracción de gas y petróleo
22	Servicios
23	Construcción
31,32,33	Manufacturas
42	Venta al por mayor
44,45	Comercio minorista
48,49	Transporte y almacenamiento
51	Información
52	Finanzas y seguros
53	Bienes raíces, alquiler y arrendamiento
54	Servicios profesionales, científicos y técnicos
55	Gestión de empresas
56	Gestión de residuos y saneamientos
61	Educación
62	Sanidad y asistencia social
71	Arte, entretenimiento y ocio
72	Servicios de alojamiento y comida
81	Otros servicios (excepto Administración Pública)
92	Administración Pública

Tabla 6.3 - Categorías generales para el NAICS

Para esta experimentación, al igual que para la segunda, se han creado las nuevas clases *SeleccionConversionExperimentacion3.java*, *RecuentoExperimentacion3.java*, y la clase *CompletacionExperimentacion3.java* que se adecuan a las nuevas condiciones de agrupación sobre las dos características elegidas; además, en este caso es necesario incluir también *SeleccionMapFunctionExperimentacion3.java* porque es la encargada de aplicar la función de transformación y en ella se seleccionan las dos primeras cifras del NAICS a transformar en un valor discreto para el fichero de salida.

Una vez generados los ficheros de salida y completadas las series con las observaciones a cero necesarias, se ha aplicado regresión múltiple lineal con los resultados mostrados en la Figura 6.3.


```

Residuals:
    Min       1Q   Median       3Q      Max
-1853    -628    -243     110   47141

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 96897.595  12474.741   7.768 8.28e-15 ***
dinero      -116.253    3.114  -37.330 < 2e-16 ***
mes          -3.395    3.677   -0.923  0.356
year        -47.199    6.202   -7.610 2.82e-14 ***
naics       -33.105    1.569  -21.100 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2105 on 28219 degrees of freedom
Multiple R-squared:  0.06122,    Adjusted R-squared:  0.06109
F-statistic: 460.1 on 4 and 28219 DF,  p-value: < 2.2e-16

```

Figura 6.3 - Regresión múltiple empleando sólo presupuesto y NAICS reducido

El coeficiente de determinación ajustado en este caso ya es de 0,06109, explicando el modelo el 6,1% de la varianza de los datos: mucho mayor que en las experimentaciones anteriores pero aún muy bajo.

Se decide que no merece la pena la pérdida de información que supone prescindir del código de agencia de contratación y la simplificación de sectores de actividad conseguidos con la reducción del NAICS para alcanzar este porcentaje aún pequeño, así que la solución final del problema va a pasar por seleccionar para cada serie de la primera experimentación el mejor modelo de aproximación en función de los resultados mostrados por la regresión lineal, los modelos ARIMA y la red de neuronas artificial.

6.4 Comparativa de la aproximación de los modelos para un conjunto de series

Como fase final de la experimentación se ha elaborado una comparativa entre los resultados de los tres modelos implementados para identificar cuál se ajusta mejor en cada caso y en qué proporción.

Para ello se ha implementado en el método *main* de la clase Java *ComparativaModelo.java* la lógica necesaria para obtener los valores para un conjunto de ternas determinadas durante el año fiscal 2015 a partir del modelo de regresión lineal, el mejor modelo ARIMA que se ajuste a esa serie y los resultados que ha ofrecido para ese terna la red de neuronas.

Se ha partido de una selección de 12.147 series (de las más de 200.000 que había en total) y para cada una de ellas se ha hecho:

- Se calcula para los 12 meses del año fiscal 2015 los valores obtenidos por regresión lineal múltiple utilizando los coeficientes de la primera experimentación. Se ha tratado de calcular el modelo desde la clase Java utilizando *JRI* para la integración con R: se cargan correctamente todos los ficheros de contratos pero el

programa finaliza de manera inesperada al calcular el modelo. Parece que es debido a un problema de memoria a la hora de almacenar el modelo final, que en R ocupa más de 400 MB. La alternativa más rápida y sencilla es partir de los coeficientes ya calculados y aplicar la ecuación, liberando a esta clase de tener que volver a recalcular el modelo. La fórmula, -redondeando a dos decimales para su representación- a aplicar es la que se detallaba anteriormente:

$$N^{\circ} \text{ Contratos} = 152.329 - 0.152 * \text{Dinero} - 0.006 * \text{Mes} - 0.078 * \text{Año} \\ + 0.002 * \text{Naics} + 0.02 * \text{Agencia}$$

- Se aplica la función *auto.arima* sobre la serie de observaciones y se obtienen los valores pronosticados para los próximos 12 meses. Esto se ha hecho completamente desde la clase Java integrándola con R e invocando la función *auto.arima* del paquete *forecast*.
- A partir del fichero de test empleado en la red de neuronas LSTM que contenía las 12.147 series para las se quería determinar sus 12 siguientes observaciones, se van leyendo desde esta clase los valores generados en el fichero de salida de test de la red.

Teniendo los datos reales de cada serie y los datos de la predicción para cada uno de los modelos, se ha aplicado la raíz cuadrada del error cuadrático medio (*RMSE*: *Root Mean Squared Error*) según se especifica en la fórmula:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Siendo n el número de observaciones -12-, y_i el valor real de la observación e \hat{y}_i el valor obtenido en la predicción del modelo aplicado.

Con esta medida, *RMSE*, es posible comparar los modelos eligiendo aquel que presente un menor error como el que mejor aproxima el total de contratos del siguiente año fiscal.

La clase *ComparativaModelo* genera un fichero de texto con los resultados obtenidos como puede apreciarse en el extracto representado en la Tabla 6.6. Por cada terna de presupuesto, NAICS y agencia se presentan los valores reales y los obtenidos en cada uno de los tres modelos, el RMSE de cada modelo y se marca con una 'X' cual de los tres presenta el error más bajo.

En el histograma de la Figura 5.9 se representaba la distribución del número de series según el número de observaciones que tenían distintas de 0 entre las 96 observaciones posibles del período fiscal 2008-2015.

De las más de 202.046 series hay 176.253 que tienen entre 1 y 20 observaciones distintas de 0. 12.874 tienen entre 21 y 40 observaciones. 7.772 tienen entre 40 y 72

observaciones. 4.017 tienen más del 75% de observaciones rellenas, entre 72 y 95 observaciones. Sólo 1.130 series tienen todas sus 96 observaciones completas. La distribución de series elegida para comparar los modelos ha sido la mostrada en la Tabla 6.4.

N.º de observaciones distintas de 0	N.º de series escogidas
Entre 1 y 20	2.000
Entre 20 y 40	2.000
Entre 40 y 72	3.000
Entre 72 y 95	4.017
Con todas las series (96)	1.130
Total	12.147

Tabla 6.4 - Cantidades de series de observaciones seleccionadas en función del número de observaciones mayores de 0

Para las series de entre 1-20, 20-40 y 40-72 se ha hecho una selección aleatoria debido a su elevado número; para aquellas que tienen igual o más de 72 observaciones se han incluido todas.

La tabla y gráfico que muestran la distribución de mejor modelo se muestra en la Tabla 6.5 y la Figura 6.4.

Modelo	Menor RMSE	Porcentaje
Regresión Múltiple	2.111	17,38%
ARIMA	5.592	46,04%
Red LSTM	4.444	36,59%

Tabla 6.5 - Resultados porcentuales de cada modelo sobre el total de series seleccionadas

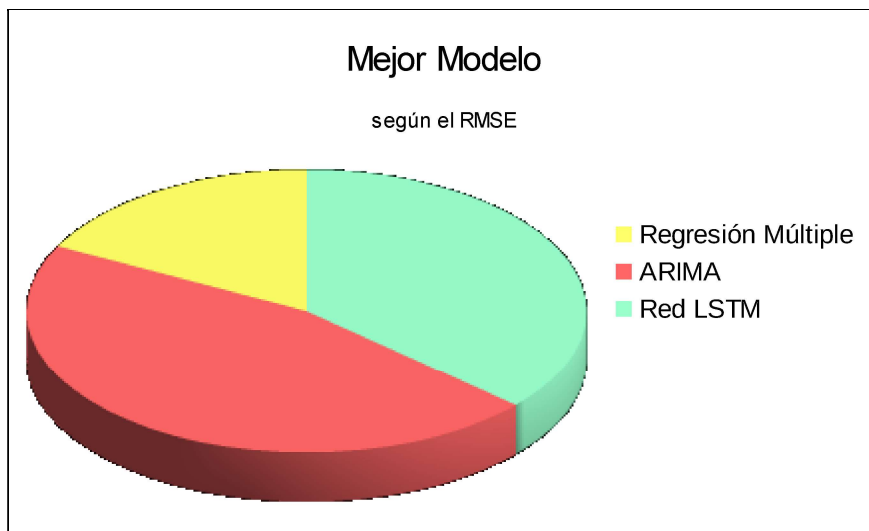


Figura 6.4 - Gráfico con el porcentaje de series que mejor cubre cada modelo

Si se calcula la media del RMSE de los mejores modelos escogidos para cada una de las series se tiene un valor de 12,98. El RMSE máximo obtenido para una serie con su mejor modelo es de 35.620,47 y el mejor es 0.

Para 97 series el RMSE es mayor de 100 y para 694 es mayor que la media calculada.

Se observa que los mayores errores se obtienen aproximando aquellas series de contratos donde cada observación es muy alta, por encima de los 1.000 contratos mensuales mientras que todos los modelos se ajustan bastante bien a las series poco pobladas, donde la mayoría de las observaciones son 0.

6.5 Matriz de validación de requisitos SW

El grado de cumplimiento de los requisitos software definidos al inicio del proyecto ha sido total para todos y cada uno de ellos. Se incluye la matriz de validación en la Figura 6.5, detallando en las observaciones en qué punto o con cuántos experimentos ha sido satisfecho el requisito.

	Ninguno	Parcial	Total	Observaciones
SR-F-01			X	
SR-F-02			X	Clase de Preselección
SR-F-03			X	Análisis Características en R
SR-F-04			X	Clase de División
SR-F-05			X	Clase de Selección y Conversión
SR-F-06			X	Clase de Completación
SR-F-07			X	3 experimentos con Regresión Múltiple
SR-F-08			X	Aplicación de la función auto.arima
SR-F-09			X	Estudio de 4 modelos de red LSTM
SR-F-10			X	Clase de Comparativa de los modelos
SR-O-01			X	
SR-O-02			X	
SR-O-03			X	
SR-O-04			X	Ficheros .properties por cada clase
SR-R-01			X	
SR-D-01			X	
SR-D-02			X	

Figura 6.5 - Matriz Grado Cumplimiento Req. SW

PRESUP.TO	NAICS	AGENCIA	TIPO	OCT14	NOV14	DIC14	ENE15	FEB15	MAR15	ABR15	MAY15	JUN15	JUL15	AGO15	SEP15	RMSE	MEJOR
3	1601	290	REAL	6.247	5.022	5.136	5.631	5.282	5.552	5.670	5.089	5.384	5.440	4.894	4.693		
			REG.MUL.	2,48	2,47	2,46	2,45	2,44	2,44	2,43	2,42	2,42	2,41	2,41	2,40	5.349,06	
			ARIMA	29,77	49,36	62,25	70,75	76,33	80,01	82,44	84,03	85,08	85,77	86,23	86,53	5.279,16	X
			RNN	18	13	14	15	16	17	18	16	18	16	18	23	5.334,72	
3	1636	290	REAL	29.508	28.378	36.036	35.111	25.171	1	1	41.312	48.249	47.641	49.395	40.682		
			REG.MUL.	2,54	2,53	2,53	2,51	2,51	2,50	2,50	2,49	2,48	2,48	2,47	2,47	35.635,06	
			ARIMA	5,81	5,81	5,81	5,81	5,81	5,81	5,81	5,81	5,81	5,81	5,81	5,81	35.632,11	
			RNN	20	14	16	16	18	19	20	19	20	18	20	25	35.620,47	X
5	1627	201	REAL	0	0	0	0	0	0	0	0	0	0	1	1		
			REG.MUL.	0,45	0,45	0,44	0,43	0,42	0,41	0,41	0,40	0,40	0,39	0,39	0,38	0,46	
			ARIMA	0,41	0,41	0,41	0,41	0,41	0,41	0,41	0,41	0,41	0,41	0,41	0,41	0,45	
			RNN	0	0	0	0	0	0	0	0	0	0	1	1	0,00	X
4	1568	284	REAL	0	0	1	0	0	0	1	1	1	0	1	1		
			REG.MUL.	2,14	2,14	2,13	2,12	2,11	2,10	2,10	2,09	2,09	2,08	2,08	2,07	1,68	X
			ARIMA	3,35	3,35	3,35	3,35	3,35	3,35	3,35	3,35	3,35	3,35	3,35	3,35	2,89	
			RNN	2	1	1	2	2	2	2	2	2	2	2	4	1,69	
3	1501	264	REAL	37	30	26	27	25	29	28	27	25	25	27	26		
			REG.MUL.	1,77	1,77	1,76	1,75	1,74	1,73	1,73	1,72	1,72	1,71	1,71	1,7	26,14	
			ARIMA	3,84	4,94	5,56	5,9	6,07	6,16	6,2	6,23	6,24	6,24	6,25	6,25	22,19	X
			RNN	5	4	4	4	4	5	5	4	5	5	6	8	23,01	
5	1296	290	REAL	0	0	20	24	39	37	7	2	0	3	0	0		
			REG.MUL.	1,6	1,6	1,59	1,58	1,57	1,57	1,56	1,55	1,55	1,54	1,54	1,53	17,19	
			ARIMA	9,27	15,98	18,84	18,13	15,14	11,46	8,43	6,77	6,59	7,49	8,88	10,18	12,49	X
			RNN	9	7	7	7	8	9	10	8	9	9	10	12	15,11	

Tabla 6.6- Extracto del fichero generado en la clase de ComparativaModelo, donde se recoge el desempeño de cada modelo por cada serie temporal seleccionada y su desviación respecto al valor real para los meses del año fiscal 2015

7. Conclusiones

Atendiendo a los resultados obtenidos durante las etapas de implementación y experimentación se van a enumerar las principales conclusiones obtenidas sobre la idoneidad de las metodologías y herramientas utilizadas para el tratamiento del problema y se va a exponer por qué los resultados obtenidos están lejos del objetivo deseado al inicio del proyecto.

Apache Flink ha resultado ser una herramienta sencilla, útil y eficaz en el tratamiento de los grandes ficheros de datos de contratos de los que se partía. Ofrece un API lo suficientemente extensa, madura y documentada como para no ser necesario un conocimiento profundo de la misma para poder abordar el problema del tratamiento de estos ficheros sin grandes dificultades. Las carencias principales encontradas han sido a la hora de realizar agregaciones, por no soportar tipos complejos de datos y no estar aún implementados algunos tipos de agregación, como la cuenta del número de registros agrupados. En cuanto a su empleo para realizar aprendizaje automático aún está falto de recorrido y otras alternativas -como Apache Spark, por ejemplo- constituyen opciones más completas en cuanto a disponibilidad de capacidades; el hecho de sólo poderse implementar en Scala también resta atractivo al uso de FlinkML como librería para *machine learning*.

El rendimiento de Apache Flink tratando los ficheros ha sido muy satisfactorio teniendo en cuenta las limitaciones de procesamiento del equipo donde se han ejecutado los experimentos, por lo que se recomendaría su uso en problemas más complejos. No obstante hay que tener especial cuidado a la hora de emplear el paralelismo, puesto que indicar un mayor índice no siempre garantiza unos mejores resultados, como ha ocurrido en la implementación ejecutando sobre varios ficheros simultáneamente y configurando un índice de paralelización de 4.

Al final el problema de predicción de contratos se ha reducido a un problema de tratamiento de series temporales. La dificultad en el tratamiento de una serie de contratación reside principalmente en que se trata de un fenómeno socioeconómico y que no siempre obedece a criterios objetivos y cuantificables para los que se puede medir su impacto en la evolución del índice de contratación. Así, el número de licitaciones se ve alterado por cambios en el ciclo económico, decisiones políticas de inversión o fenómenos naturales o geopolíticos que obligan a modificar las prioridades de gasto; todo esto añade incertidumbre a la serie y dificulta su modelización contando sólo con la información de los contratos.

Los resultados pronosticados por los tres modelos testados están lejos de ser lo suficientemente aceptables como para utilizarse en un modelo productivo del que obtener un rendimiento económico: es difícil convencer a alguien para que invierta su dinero en una predicción tan poco exacta.

Los tres modelos funcionan razonablemente bien para series cuyas observaciones son en la mayoría de instantes de tiempo de 0 contratos. Estas series, según la agrupación de características escogidas, son mayoría -el 87%-, pero también son las que menos interés representan: si casi nunca hay contratos de un tipo lo más probable que en el futuro

tampoco los haya. El rendimiento de los tres modelos se empobrece para aquellas series que tienen valores mayores de 0 en más del 75% de observaciones, y especialmente para aquellas que tienen muchos contratos mensuales: se está penalizando esta minoría de series impactadas por la mayoría de series con observaciones a 0. Es por esto que aunque el error medio es de menos de 13 contratos anuales se dispara hasta valores mucho más altos, entre 35.000 y 100 contratos anuales, para casi una centena de series en total.

Revisando el desempeño individual de los modelos se aprecia que la regresión múltiple lineal es la que peor comportamiento presenta mientras que el tratamiento de las series temporales de manera individual con ARIMA ofrece el mejor resultado en casi la mitad de los casos. La superioridad de ARIMA sugiere que los valores de la serie no se ven alterados por características externas sino por el propio fenómeno en sí: es decir, que un valor en un momento determinado se ve principalmente condicionado por los valores anteriores de la serie sin atender a otros factores exógenos a esta. También es cierto que con ARIMA se está tratando cada serie temporal de una terna de {presupuesto, NAICS, agencia} de manera individual e independiente del resto, por lo que cada serie es ajena a las otras y no se ve condicionada por los valores que aquellas presenten.

La red de neuronas LSTM aplicada a este problema presenta un comportamiento aceptable para series con la mayoría de observaciones a 0 pero también lo hace para aquellas que tienen la mayoría de observaciones distintas de 0 con una horquilla de contratos de 5-50 mensuales; no consigue aproximar tan bien cuando se trata de observaciones con miles de contratos mensuales. A priori es el modelo que más potencial de mejora presenta, pues siendo una estructura de red sencilla y con un entrenamiento rápido –alrededor de una hora- es capaz de generalizar casi tan bien como ARIMA, sin necesidad de individualizar cada serie.

El uso de R como lenguaje para la implementación de los modelos ARIMA y su posibilidad –y facilidad- de integración con Java, conjuntamente con la existencia del paquete *forecast*, ha facilitado enormemente la implementación masiva del modelo ARIMA empleando la función *auto.arima*, que evita la tarea previa de tener que descomponer la serie en su parte estacionaria y no estacionaria y después aplicar el modelo. Así, una tarea que podría haber sido tediosa de tenerse que identificar y calcular previamente las componentes para luego aplicar el modelo, se ha convertido en un proceso automático que además requiere pocas líneas de código para su implementación.

La implementación de una red de neuronas con Keras y TensorFlow ha sido satisfactoria especialmente por dos factores: la facilidad de instalación y configuración de las herramientas y la sencillez del API de Keras para diseñar un modelo y hacerlo correr sobre TensorFlow. Para problemas con un poco más de complejidad, en función de aspecto del que se trate, se recomiendan dos alternativas:

- Si la complejidad viene dada por el número de elementos del conjunto de entrenamiento, se recomienda emplear utilizar el entrenamiento de la red en modo *batch* que facilita Keras con los métodos *model.train_on_batch* y *model.test_on_batch*.
- Si la complejidad viene porque la red de neuronas empleadas utiliza un mayor número que las capas de este problema –o con el mismo número de capas pero más neuronas en cada una de ellas-, se recomienda dimensionar adecuadamente la memoria disponible para Docker o para la máquina virtual, –si es que se ha descartado su instalación en un entorno Unix-, para no sufrir continuos problemas de memoria con TensorFlow.

8. Líneas Futuras

Contando con los ficheros de contratación como materia prima y las clases implementadas se pueden hacer muchas actividades complementarias a las ya realizadas en este proyecto y abordar otras que han quedado fuera del alcance del mismo.

Se pueden realizar muchas más experimentaciones seleccionando otros conjuntos de características con las que identificar los contratos y agrupar por ellas. Aunque es posible que aplicando regresión no mejore significativamente la calidad de los modelos obtenidos, una vez construidos los nuevos ficheros aptos para aplicarles ARIMA o utilizarlos para entrenar una red neuronal se podrían obtener mejores resultados que hasta ahora. Es más, experimentar con otros modelos de tratamiento de series temporales o con alternativas de los ya estudiados puede derivar en modelos que presenten un menor margen de error respecto a la serie real. Por ejemplo, implementar otros diseños que utilicen otras estructuras y configuración para la red de neuronas artificial, o plantear alguno modelo híbrido entre ARIMA y una red de neuronas como en [69].

Es importante no olvidar que la manera en la que se han agrupado los contratos ya está implicando un sesgo a los modelos en tanto que se está ajustando la realidad de los datos a un modelo más simple y manejable. Por ejemplo, si se hubiesen planteado rangos de presupuesto con horquillas de miles de dólares más amplias habría menos series de observaciones y éstas a su vez tendrían menos mediciones a cero, puesto que los contratos se habrían agrupado más al haber menos opciones disponibles. Se aprecia que los modelos ajustan mejor las series con más observaciones a cero que aquellas series que han mostrado contrataciones todos los meses; por esto, se propone una alternativa de experimentación: entrenar los modelos presentados pero de manera independiente con dos conjuntos de datos diferenciados: por un lado las series con más del 75% de sus observaciones con valores mayores de cero y por otro lado el resto. Esto sería interesante para comparar los modelos resultantes y su capacidad de predicción en el caso de la regresión múltiple y la red de neuronas, viendo si teniendo dos modelos separados según las características de la serie se ajusta mejor que habiendo generalizado con todas (en el caso de ARIMA no tiene sentido, pues se ha aplicado a cada serie temporal de manera independiente).

Otra línea de trabajo adicional: no se ha comprobado dentro de este proyecto cómo actúan los modelos seleccionados en cada serie en real, por lo que se podría comprobar fácilmente como ajusta cada uno de ellos los datos de los meses del año fiscal 2016 en curso y ver cuánto se ha desviado de la realidad de los datos.

Durante la fase de división se ha dejado de lado mucha información: la de todas aquellas líneas de los ficheros de entrada que hacían referencia a modificaciones de contratos. Se pueden aplicar también sobre éstas técnicas de aprendizaje automático para tratar de adivinar cómo va a evolucionar un contrato en curso.

Otra alternativa a la predicción aplicada en el proyecto consistiría en tratar de determinar, en base a las características de un contrato, qué proveedor de los habituales tendría más posibilidades de ganar su adjudicación.

Otros países como Australia, Grecia o Reino Unido, también tienen a disposición de terceros los datos de contratación de su administración pública. Se podría trabajar con esas fuentes y tratar de resolver el mismo problema.

Aunque por los resultados obtenidos gracias a los modelos ARIMA podría parece que la serie de observaciones sólo depende de sí misma, se podrían plantear modelos que estudiaran la dependencia entre la serie de contratación y otros fenómenos más o menos externos a ella. Por ejemplo, puesto que los presupuestos para un año fiscal se aprueban antes del comienzo de éste, se puede estudiar la influencia de un aumento o reducción del presupuesto sobre la serie de observaciones, viendo cómo afecta al número de contratos. Otra alternativa: estudiar si el cambio de gobierno afecta a la contratación, determinando si hay impacto en la serie de contratos tras una alternancia en la administración norteamericana entre Demócratas o Republicanos, y viceversa.

Atendiendo a la implementación de las clases Java para el tratamiento de los ficheros de datos y conociendo ya la secuencia y funcionalidad de cada una de ellas sería muy beneficioso *refactorizarlas* pensando en su aprovechamiento futuro para una mejora en los tiempos del proceso. Ahora mismo, cada fase está identificada con una clase distinta y prácticamente cada una utiliza como entrada los ficheros de salida generados en la clase anterior. La escritura a disco de esos ficheros consume recursos de espacio y sobre todo de tiempo de escritura teniendo en cuenta que luego esos mismos datos se van a cargar de nuevo para someterlos a nuevas transformaciones hasta llegar al resultado final. Si sólo se vuelca a fichero al finalizar el proceso y mientras tanto se trabaja con los *datasets* de Flink que los mantienen en memoria para operar con ellos – salvo las *serializaciones* a disco necesarias derivadas de las limitaciones de memoria- se ahorraría toda esa escritura intermedia a ficheros que son reutilizados posteriormente. Así, en una misma clase podría abordarse el proceso completo de transformación de los ficheros de entrada empleando para ello mucho menos tiempo.

Si se buscara una explotación comercial transformando este proyecto en un producto sería necesario implementarlo como un proceso continuo y en constante actualización, que se descargara periódicamente los nuevos ficheros de contratos, les aplicara el tratamiento y transformación correspondiente y con ellos se actualizaran los modelos de predicción incluyendo esta nueva información. Así, se podría realizar periódicamente una predicción de la evolución de la contratación y ofrecérsela a autónomos y empresas interesados en conocer cuáles y cuándo pueden ser sus próximas oportunidades de contratación con la administración pública. Esto podría capitalizarse a través de un servicio de suscripción, donde en base a unos sectores de actividad económica de interés para una empresa el sistema actúe como recomendador, enviando predicciones periódicas sobre las oportunidades de contratación pronosticadas.

9. Gestión del Proyecto

9.1 Planificación

En este apartado se hace un desglose de todas las tareas y subtarear aboradas durante la ejecución de este proyecto y se acompaña de los diagramas de Gantt que mejor ilustran la evolución del proyecto.

No se ha seguido ningún estándar para la elaboración de la planificación del proyecto y su gestión por no ser este el objetivo principal de proyecto y dado el alcance limitado del mismo, pero sí ha servido el estándar IEEE 16326:2009 [15] como aproximación en el transcurso del mismo.

Este proyecto ha tenido un enfoque particular en tanto que se conocía el origen del proyecto, -los ficheros de datos de contratación de EE. UU.-, y el destino -predecir cuándo se concederán algunos contratos del siguiente año fiscal-, pero no sé había determinado cuál sería la mejor manera de abordar el problema, ni los medios y plazos necesarios hasta conseguirlo. A medida que se han ido conociendo más detalles sobre los contratos y su información disponible se han ido aportando formas de afrontar la solución del problema dimensionando los recursos humanos, técnicos y temporales para resolverlos.

En el primer diagrama de Gantt inicial, Figura 9.1, se incluyen las tareas previstas al inicio del proyecto enfocadas al análisis de la información de contratos y de Apache Flink como herramienta de tratamiento de datos masivos y se apuntaba en líneas generales el resto. En un segundo diagrama, Figura 9.3, se incluye la planificación final con todas las tareas desarrolladas.

En cuanto a los recursos humanos empleados y su dedicación al proyecto hay que señalar que aunque es la misma persona –el autor de este proyecto– la que ha realizado todas las tareas del mismo (a excepción, obviamente, de las tareas de supervisión ejecutadas por el tutor) sí se han querido diferenciar cuatro roles distintos dependiendo de las tareas realizadas en cada momento, pues, de poder hacerse el proyecto entre varias personas diferenciadas, cada uno de estos roles podría ser asumido por un analista, programador, o analista programador independiente del resto de personas que trabajaran en la ejecución del mismo. Los roles incluidos son:

- **fparis-Documentacion:** lleva toda la adquisición de nuevo conocimiento necesaria para la realización del proyecto. El estudio y aprendizaje de los lenguajes de programación no conocidos hasta el momento, la documentación de las nuevas herramientas, la búsqueda de referencias de utilidad y la familiarización con el proceso de contratación. También es el encargado de generar toda la documentación resultado de la elaboración de este proyecto y de participar en las tareas de seguimiento del mismo. Supone un rol de Analista o Documentalista según el caso concreto.

- **fparis-Flink:** encargado de la implementación en Apache Flink de la solución necesaria para tratar los ficheros de contratos de entrada hasta convertirlos en ficheros válidos para ser evaluados por los modelos de aprendizaje y predicción elegidos. Desempeña un rol de: Programador o Analista Programador, según la fase.
- **fparis-R:** encargado de realizar el análisis de características de los contratos y la implementación en Java (*rJava*) de los modelos ARIMA y regresión múltiple. Rol de Analista Programador.
- **fparis-TensorFlow:** toda la instalación del entorno y la implementación de las redes de neuronas. Rol de Analista Programador.

□ Planificación Inicial

Como se indicaba antes, al inicio del proyecto no se tenían claras las fases del mismo, así que la planificación inicial del proyecto sólo contiene como definitiva la fase de Análisis Inicial que es cuándo se elaboró esta planificación y como estimadas el resto de tareas.

1. Análisis Inicial

Engloba todas las tareas que se llevaron a cabo en la toma de contacto con el dominio del problema y con la herramienta de Apache Flink.

1.1 Revisión Información Contratación

Lectura de toda la documentación relacionada con el procedimiento, información y formatos de la información de contratación que está disponible en la web del gobierno norteamericano.

1.2 Reunión de Kick-Off

Reunión con el tutor para el arranque del proyecto, definiendo el alcance del mismo, recursos de partida y plazos a cumplir.

1.3 Introducción a Apache Flink

Estudio de la documentación de Apache Flink y aplicación de la guía rápida para construir un ejemplo funcional. Familiarización con la herramienta.

1.4 Prueba de concepto con Flink (PoC, *Proof of Concept*)

Construcción de una clase Java que partiendo de uno de los ficheros de contratos fuera capaz de generar un fichero de salida aplicando algunas transformaciones a los datos.

1.5 Aprendizaje Automático con Apache Flink

Revisión del estado de la librería de FlinkML y estudio de las posibilidades que ofrece en este momento y conociendo sus limitaciones.

2. Tratamiento de los ficheros de Contratos

Se estimaron 14 días laborables para esta fase, que se descomponía en las tareas de: análisis y selección de las características, implementación de las transformaciones y documentación del proceso. En estas subtarefas se engloba el desarrollo de las transformaciones necesarias sobre los datos de entrada para generar unos ficheros válidos para la predicción, y la documentación de todo el proceso.

3. Predicción

Se estimaron 27 días laborables para esta fase, que incluía las tareas de: selección de modelos de predicción, experimentación, análisis de resultados y conclusiones y documentación del proceso. En esta fase se contemplaba la búsqueda y selección de modelos de predicción sobre los datos de contratación y el análisis comparativo de los resultados obtenidos sobre ellos. Aún no se habían estudiado y decidido los tres modelos que finalmente constan en la experimentación de este proyecto.

4. Documentación

Para esta fase se estimaron 18 días laborables en los que debería elaborarse la documentación que faltara para esta memoria y la presentación de la misma.

En Figura 9.1 se incluye el desglose de tareas y subtarefas el diagrama de Gantt de la planificación inicial.

□ Planificación Final

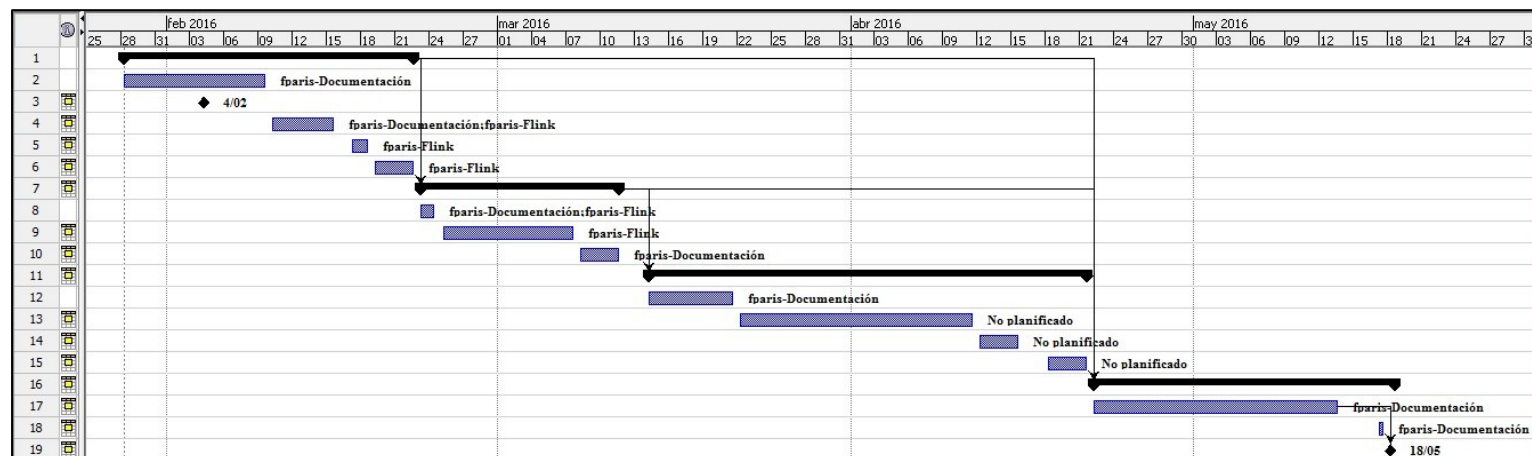
Las subtarefas del Análisis Inicial ya eran definitivas en la planificación inicial.

La comparativa entre la planificación inicial y la final puede desglosarse como:

- En la planificación final la tarea de “Tratamiento de los ficheros de Contratos” se ha dividido finalmente en dos grandes tareas: “Tratamiento de los ficheros de Contratos” y “Generación ficheros para Predicción”. Inicialmente estaban previstos 14 días y se han consumido un total de 22 días, principalmente en el estudio de la información de los contratos para determinar las mejores características a emplear en la predicción y la transformación necesaria.
- La tarea de “Predicción” de la planificación de partida se ha desdoblado en “Estudio de alternativas para la predicción”, “Experimentación con Regresión Múltiple”, “Experimentación con Modelos ARIMA”, “Experimentación con Redes Neuronales LSTM” y “Análisis de los resultados obtenidos”. Se habían estimado 27 días y se han empleado un total de 36 días.

	Nombre	Duración	Inicio	Terminado	Nombres del Recurso	Predecesores
1	Análisis inicial	18 days	28/01/16 8:00	22/02/16 17:00		
2	Revisión información contratación	9 days	28/01/16 8:00	9/02/16 17:00	fparis-Documentación	
3	Reunión Kick-off	1 day	4/02/16 8:00	4/02/16 17:00	fparis-Documentación	
4	Introducción a Apache Flink, Documentación	4 days	10/02/16 8:00	15/02/16 17:00	fparis-Documentación; fparis-Flink	
5	Prueba de concepto con Flink	2 days	17/02/16 8:00	18/02/16 17:00	fparis-Flink	
6	Aprendizaje Automático con Apache Flink, Documentación	2 days	19/02/16 8:00	22/02/16 17:00	fparis-Flink	
7	Tratamiento de los ficheros de Contratos	14 days?	23/02/16 8:00	11/03/16 17:00		1
8	Análisis y selección de características	1,5 days?	23/02/16 8:00	24/02/16 13:00	fparis-Documentación; fparis-Flink	
9	Implementación de las transformaciones	8 days?	25/02/16 8:00	7/03/16 17:00	fparis-Flink	
10	Documentación del proceso	4 days?	8/03/16 8:00	11/03/16 17:00	fparis-Documentación	
11	Predicción	27 days?	14/03/16 8:00	21/04/16 17:00		7
12	Selección de modelos	6 days?	14/03/16 8:00	21/03/16 17:00	fparis-Documentación	
13	Experimentación	13 days?	22/03/16 8:00	11/04/16 17:00	No planificado	
14	Análisis de resultados y conclusiones	4 days?	12/04/16 7:00	15/04/16 17:00	No planificado	
15	Documentación del proceso	4 days?	18/04/16 7:00	21/04/16 17:00	No planificado	
16	Documentación	18 days?	22/04/16 8:00	18/05/16 17:00		1;7;11
17	Memoria	16 days?	22/04/16 8:00	13/05/16 17:00	fparis-Documentación	
18	Presentación	1 day?	16/05/16 8:00	17/05/16 17:00	fparis-Documentación	
19	Entrega	1 day?	18/05/16 8:00	18/05/16 17:00		17

Figura 9.1 - Desglose de tareas y diagrama de Gantt de la fase de planificación Inicial



- La tarea de Documentación se mantiene como se planificó, aunque se han consumido sólo 12 de los 18 días planificados. El ahorro se debe a que parte de la documentación que se había contemplado en la planificación inicial se ha ido incluyendo en la memoria del proyecto en tareas anteriores, sin esperar a finalizar la implementación y experimentación para la escritura del documento.

El desglose de tareas y subtareas de las nuevas definitivas incluidas en esta planificación final es la que sigue, considerando como primera tarea la de Análisis Inicial ya descrita:

2. Tratamiento de los ficheros de contratos

2.1 Selección de características representativas

Se ha realizado en este punto la revisión de las 250 características de partida y su posible utilidad. Para entender mejor muchas de ellas ha sido necesario llevar a cabo un proceso de inducción a partir de la información de conjuntos de contratos para ver qué clase de valores tomaban y cuál podía ser su significado y la relación con el objeto del contrato u otras características de este.

2.2 Análisis estructura y formato de los ficheros

Se ha incluido en esta subtask el tratamiento del fichero de partida para dividir y seleccionar las características necesarias y la realización en Java de la clase de Preselección. En este punto es donde se ha descubierto que no todas las líneas del fichero están bien formadas y que la estructura de encerrar cada campo entre dobles comillas no siempre facilita su tratamiento.

2.3 Estudio de la distribución y naturaleza de las características

Con la extracción masiva de los valores de las que se han analizado y su estudio estadístico llevado a cabo en R se ha identificado el dominio donde toman valor y el índice de repetición para cada una de ellas, así como servir de base para su codificación en valores discretos aptos para un modelado matemático.

2.4 Implementación transformación de los ficheros

En este punto se ha incluido la elaboración de las clases Java que se han encargado de realizar el tratamiento de datos a partir de los ficheros con las características preseleccionadas. Las implementaciones de las clases de División, y Selección-Conversión se engloban en este punto.

3. Generación Ficheros para Predicción

3.1 Consolidación de series temporales completas

Las fases de Recuento y Completación no se habían contemplado hasta no comprobar que faltaban observaciones en las series para los meses en los que no se habían concedido contratos para una determinada terna {presupuesto, agencia, NAICS}.

En esta subtarea se implementan las clases de Recuento y Completación y se analiza la distribución del número de observaciones mayores de 0 en todas las series.

3.2 Documentación del proceso de tratamiento y generación

En esta subtarea se redacta todo el proceso relacionado con las tareas elaboradas utilizando Apache Flink que se emplea en el apartado de Implementación de esta memoria.

4. Estudio alternativas para predicción

Se incluyen los tres métodos con lo que se va a tratar de aproximar la secuencia de concesión de contratos. Otras alternativas, como el uso de la librería de FlinkML se descartaron antes de abordar este proceso.

4.1 Regresión lineal con múltiples predictores

En esta subtarea se analiza la idoneidad de aplicar un modelo de regresión múltiple sobre los contratos. Se determina que la manera más sencilla de aplicarlo es desde R y se recopilan y estudian los principales parámetros y gráficos a estudiar del modelo resultante para evaluar su calidad.

Se ejecuta la primera prueba sobre los datos agrupados por presupuesto, agencia y NAICS.

4.2 Series temporales con Modelo ARIMA

Se decide aplicar este modelo para determinar si ajusta mejor la serie de contratos que la regresión lineal y determinar si las observaciones son más dependientes entre sí que con las características de los contratos. Se estudia el fundamento teórico del modelo, se hace alguna prueba en R con alguna serie determinada y se decide aplicar el modelo de manera automática con la función *auto.arima* que evita tener que calcular las componentes estacionarias y no-estacionarias de manera previa al modelado.

4.3 Redes Neuronales LSTM con TensorFlow

Existía previamente un interés por probar la herramienta de Google para aprendizaje automático y se aprovechó esta subtarea para analizarla y descubrir que con redes de neuronas LSTM se habían abordado con éxito problemas de predicción de

series temporales. Se estudia también Keras y se decide emplearlo por resultar más sencillo que TensorFlow.

4.4 Documentación del Estado del Arte de estas metodologías

Toda la información recopilada sobre los distintos métodos y herramientas abordados en las subtarefas previas se documenta para incluirlas en el apartado del Estado del Arte de esta memoria.

5. Experimentación con Regresión Múltiple

5.1 Pruebas con combinaciones de múltiples predictores

Ante la pobre calidad del modelo obtenido a partir de los datos agrupados por presupuesto, agencia y NAICS se decide estudiar la aproximación mediante otros dos modelos: agrupando sólo por presupuesto y NAICS (Experimentación2) y reduciendo el NAICS a los sectores generales y presupuesto (Experimentación3).

5.2 Documentación de resultados

Se analizan los modelos obtenidos y se documentan los resultados de la regresión y las clases utilizadas en la experimentación.

6. Experimentación con Modelos ARIMA

6.1 Pruebas con `auto.arima`

Se estudia cómo funciona la función *auto.arima* y se aplica a algunas series de las obtenidas tras la completación. Se analiza cómo aplicarlo de manera automática a un conjunto masivo de series determinando que la mejor manera es hacerlo utilizando R desde una aplicación Java.

6.2 Implementación en R desde Java

Se instala *rJava* y se implementa la ejecución de la función *auto.arima* desde una clase Java.

6.3 Documentación de resultados

Se documenta el uso de la función *auto.arima* y el desarrollo de la clase Java en los apartados de Implementación y Experimentación de esta memoria.

7. Experimentación con redes neuronales LSTM

7.1 Instalación y configuración de Keras y TensorFlow

Se instala Docker como entorno donde ejecutar Keras sobre TensorFlow dada la facilidad que presenta para su uso en un sistema operativo Windows. Se prueba alguno de los ejemplos que se incluyen con la distribución de Keras que se deja configurado para que utilice TensorFlow como *backend*.

7.2 Prueba de Concepto

Se prueban algunos ejemplos de casos de uso que utilizan series temporales y redes LSTM para comprobar cómo funciona y los resultados que se obtienen.

7.3 Entrenamiento y test de redes neuronales

Se configura la red neuronal específica para modelar el dominio de los contratos y se evalúa el resultado de las variantes probadas eligiendo entre ellas el modelo definitivo.

7.4 Documentación de resultados

Se incluye en la Implementación y Experimentación todas las acciones acometidas con las redes de neuronas y las herramientas empleadas, incluyendo los resultados obtenidos.

8. Análisis de los resultados obtenidos

8.1 Selección modelo más adecuado para cada conjunto

En esta subtarea se recoge el esfuerzo invertido en la selección de un subconjunto de series de las más de doscientas mil iniciales para estudiar su aproximación con los tres métodos abordados, y la implementación de la clase Java encargada de ello, mediante la aplicación del modelo de regresión múltiple, la función *auto.arima* y los resultados obtenidos por la red de neuronas.

8.2 Documentación de resultados

Se recopilan, analizan y representan los datos obtenidos en esta última fase de la experimentación.

9. Documentación

9.1 Preparación del resto de la memoria de trabajo

Se complementan los apartados de Estado del Arte, Implementación y Experimentación con más información y gráficos. Se elabora el resto de apartados que aún no se habían acometido hasta completar la totalidad del documento.

9.2 Elaboración de la presentación para la defensa

Se realiza la presentación que se empleará en la defensa de este proyecto.

9.2 Presupuesto

Se incluye del detalle de los costes incurridos durante la elaboración de este proyecto, desglosándolos en partidas de recursos humanos y de hardware y licencias software.

□ Costes de personal

El número de horas invertidas diariamente ha sido de una media de 6, por lo que por cada una de las fases del proyecto se tiene el total de horas de la Tabla 9.1.

Se han agrupado las tareas del Gantt en cuatro grandes fases. Se ha estimado una dedicación por parte del Jefe del Proyecto en tareas de dirección y supervisión del 10% total de la dedicación invertida en el resto de tareas.

Fase	Días	Horas Empleadas
Análisis (Inicial y Estudio alternativas)	31	186
Implementación (Tratamiento y Generación ficheros)	22	132
Experimentación y análisis resultados	23	138
Documentación (Memoria, Presentación)	12	72
Jefatura del Proyecto (10%)	8,8	52,8
Total	96,8	580,8

Tabla 9.1 - Número de horas dedicadas a cada fase

Atendiendo a los roles que han participado en cada fase y a los costes por hora de trabajo estimados se tiene la Tabla 9.2 donde se calcula el coste total en recursos humanos del proyecto.


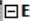


		Nombre	Duración	Inicio	Terminado	Nombres del Recurso	Predecesores
1		 Análisis inicial	18 days	28/01/16 8:00	22/02/16 17:00		
2		Revisión información contratación	9 days	28/01/16 8:00	9/02/16 17:00	fparis-Documentación	
3		Reunión Kick-off	1 day	4/02/16 8:00	4/02/16 17:00	fparis-Documentación	
4		Aproximación a Apache Flink: Documentación, Instalación	4 days	10/02/16 8:00	15/02/16 17:00	fparis-Documentación	
5		Prueba de concepto con Flink	2 days	17/02/16 8:00	18/02/16 17:00	fparis-Flink	
6		Aprendizaje Automático con Apache Flink, Documentación	2 days	19/02/16 8:00	22/02/16 17:00	fparis-Documentación	
7		 Tratamiento de ficheros de contratos	15 days	23/02/16 8:00	14/03/16 17:00		1
8		Selección de características representativas	4 days	23/02/16 8:00	26/02/16 17:00	fparis-Documentación	
9		Análisis estructura y formato de los ficheros	3 days	1/03/16 8:00	3/03/16 17:00	fparis-Documentación	
10		Estudio de la distribución y naturaleza de las características	4 days	4/03/16 8:00	9/03/16 17:00	fparis-R	8
11		Implementación transformación de los ficheros	3 days	10/03/16 8:00	14/03/16 17:00	fparis-Flink	8;9;10
12		 Generación ficheros para predicción	7 days	15/03/16 8:00	23/03/16 17:00		7
13		Consolidación de series temporales completas	2 days	15/03/16 8:00	16/03/16 17:00	fparis-Flink	
14		Documentación del proceso de tratamiento y generación	5 days	17/03/16 8:00	23/03/16 17:00	fparis-Documentación	
15		 Estudio alternativas para predicción	13 days	29/03/16 8:00	14/04/16 17:00		12
16		Regresión lineal con múltiples predictores	2 days	29/03/16 8:00	30/03/16 17:00	fparis-Documentación;fparis-R	
17		Series temporales con Modelo ARIMA	3 days	31/03/16 8:00	4/04/16 17:00	fparis-Documentación;fparis-R	
18		Redes Neuronales LSTM con TensorFlow	4 days	5/04/16 8:00	8/04/16 17:00	fparis-Documentación;fparis-TensorFlow	
19		Documentación del Estado del Arte de estas metodologías	4 days	11/04/16 8:00	14/04/16 17:00	fparis-Documentación	16;17;18
20		 Experimentación con Regresión Múltiple	3 days	15/04/16 8:00	19/04/16 17:00		
21		Pruebas con combinaciones de múltiples predictores	2 days	15/04/16 8:00	18/04/16 17:00	fparis-Flink;fparis-R	
22		Documentación de resultados	1 day	19/04/16 8:00	19/04/16 17:00	fparis-Documentación	21
23		 Experimentación con Modelos ARIMA	5 days	20/04/16 8:00	26/04/16 17:00		
24		Pruebas con auto.arima	2 days	20/04/16 8:00	21/04/16 17:00	fparis-R	
25		Implementación en R desde Java	1 day	22/04/16 8:00	22/04/16 17:00	fparis-R	
26		Documentación de resultados	2 days	25/04/16 8:00	26/04/16 17:00	fparis-Documentación	25
27		 Experimentación con Redes Neuronales LSTM	11 days	27/04/16 8:00	12/05/16 17:00		
28		Instalación y configuración de Keras y TensorFlow	3 days	27/04/16 8:00	29/04/16 17:00	fparis-TensorFlow	
29		Prueba de concepto	2 days	3/05/16 8:00	4/05/16 17:00	fparis-Documentación;fparis-TensorFlow	
30		Entrenamientos y test de redes neuronales	4 days	5/05/16 8:00	10/05/16 17:00	fparis-TensorFlow	
31		Documentación de resultados	2 days	11/05/16 8:00	12/05/16 17:00	fparis-Documentación	30
32		 Análisis de los resultados obtenidos	4 days	13/05/16 8:00	18/05/16 17:00		20;23;27
33		Selección modelo más adecuado para cada conjunto	3 days	13/05/16 8:00	17/05/16 17:00	fparis-Flink;fparis-R;fparis-TensorFlow	
34		Documentación de resultados	1 day	18/05/16 8:00	18/05/16 17:00	fparis-Documentación	33
35		 Documentación	12 days	19/05/16 8:00	3/06/16 17:00		20;23;27;32
36		Preparación del resto de la memoria de trabajo	10 days	19/05/16 8:00	1/06/16 17:00	fparis-Documentación	
37		Elaboración de la presentación para la defensa	2 days	2/06/16 8:00	3/06/16 17:00	fparis-Documentación	

Figura 9.2 - Desglose de tareas de la Planificación Final

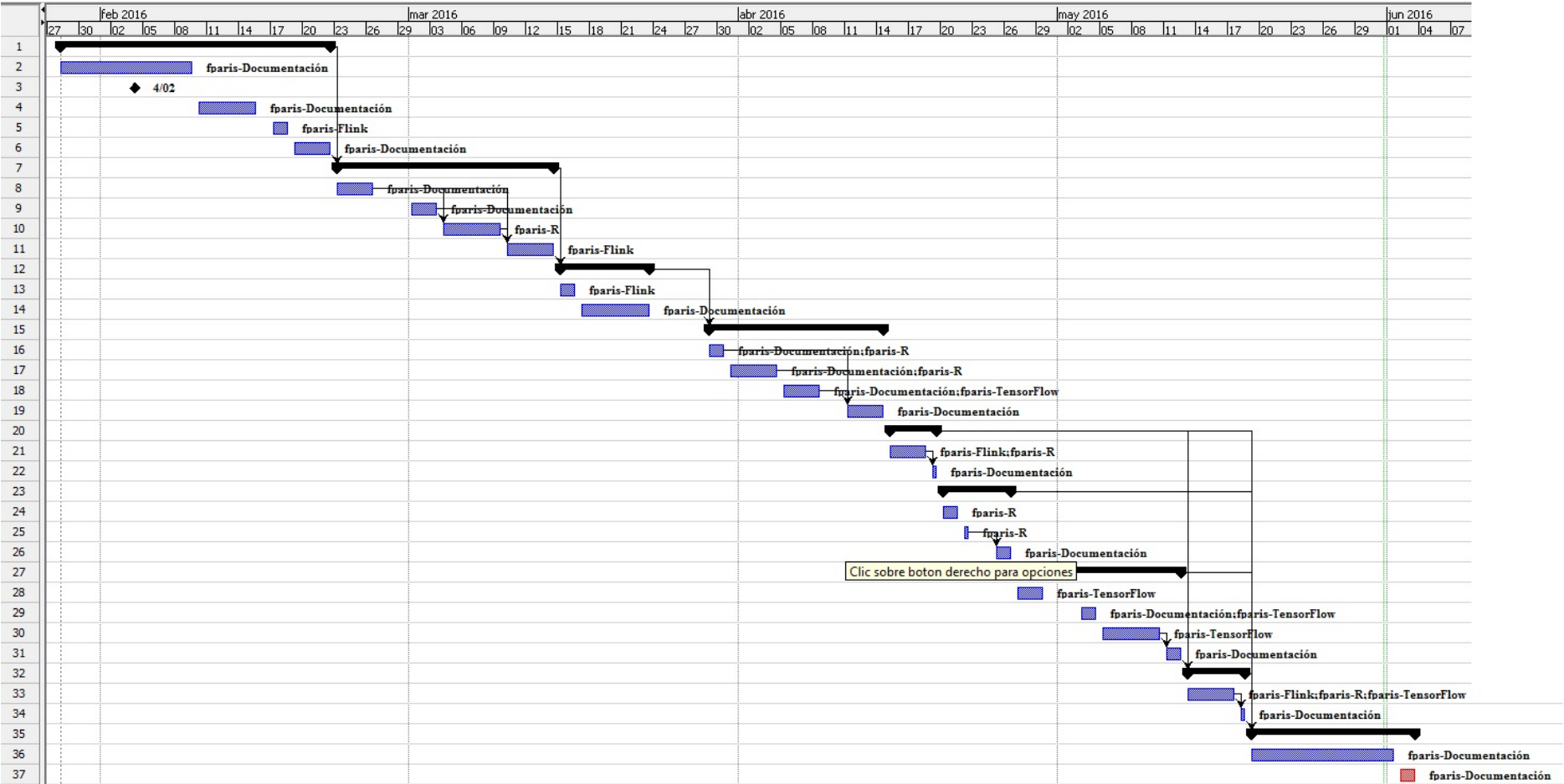


Figura 9.3 - Diagrama de Gantt con las tareas de Planificación final

Fase	Rol	Horas Empleadas	Coste €/Hora	Total
Análisis (Inicial y Estudio alternativas)	Analista	186	31	5.766 €
Implementación (Tratamiento y Generación ficheros)	Programador	132	27	3.564 €
Experimentación y análisis resultados	Analista Programador	138	31	4.278 €
Documentación (Memoria, Presentación)	Documentalista	72	24	1.728 €
Jefatura del Proyecto (10%)	JP/Tutor	52,8	45	2.376 €
			Total	17.712 €

Tabla 9.2 - Costes de recursos humanos

□ Costes de hardware y software

Para la elaboración de todas las tareas relacionadas con este proyecto se ha utilizado un portátil *Toshiba Satellite Pro C70-B-34T* con sistema operativo Windows 10 preinstalado, por lo que no hay coste por separado asociado a la licencia de éste. El coste de este ordenador asciende a 684,54 € sin incluir impuestos.

Descripción	Coste	% Dedicado	Meses	Período de Depreciación	Coste imputable
Toshiba Satellite Pro C70-B-34T	684,54 €	100	4	60	45,64 €

Tabla 9.3 - Costes de hardware

Los distintos programas empleados se desglosan en la tabla Tabla 9.4 incluyendo su modo de licenciamiento y coste aplicado.

Todas las aplicaciones empleadas son de software libre o gratuitas para el uso específico que se les da en este proyecto, excepto la aplicación ofimática de Microsoft que tiene un coste de 230,58€ antes de impuestos.

Existe una edición de pago de R Studio para empresas pero en este caso se ha empleado la versión gratuita.

Software	Modo Licenciamiento	Coste asociado
Microsoft Office Hogar y Empresas 2016	Pago Único	230,58 €
ProjectLibre	Common Public Attribution License	0,00 €
Eclipse Mars	Eclipse Public License	0,00 €
R	GNU General Public License	0,00 €
R Studio	AGPL v3	0,00 €
Docker	Apache 2.0	0,00 €
Keras	MIT License (MIT)	0,00 €
TensorFlow	Apache 2.0	0,00 €
Java JDK	Con derechos de autor [16]	0,00 €
Python	Compatible GPL	0,00 €
Apache Flink	Apache 2.0	0,00 €
	Total	230,58 €

Tabla 9.4 - Costes de software y licencias

□ Resumen de Costes

Agrupando los costes identificados de recursos humanos y materiales e incluyendo unos costes indirectos del 20%, añadiendo además los impuestos correspondientes se obtiene el total del montante económico del proyecto, mostrado en la Tabla 9.5.

Descripción	Cantidad
Recursos humanos	17.712,00 €
Costes HW+SW	276,22 €
Costes indirectos	3.597,64 €
Subtotal	21.585,86 €
Impuestos:	
IVA 21%	4.533,03 €
Total	26.118,89 €

Tabla 9.5 - Resumen presupuesto total

10. Anexo I: Manual de Usuario

Para hacer un uso correcto de cada una de las clases que resuelven las distintas fases a las que se han sometido los ficheros de contratos de entrada es necesario conocer los parámetros de ejecución de las mismas recogidos en cada caso en un fichero de propiedades independientes.

Además de esto, para configurar correctamente el proyecto no hay que olvidar importar las librerías de Apache Flink.

□ Clase **Preseleccion.java**

- ***pathFicheroCabecera***: ruta completa al fichero que contiene una única línea: la línea de cabecera de un fichero de contratos.
- ***pathFicheroCaracteristicas***: ruta completa al fichero que contiene el listado de características que hay que seleccionar de entre las que figuran en la cabecera. Cada característica se escribe en una línea independiente, sin incluirla entre comillas dobles.
- ***pathDataFeeds***: ruta completa al directorio de ficheros .csv con los datos originales de los contratos. Como todos los ficheros tienen el mismo nombre pero se diferencia en el año, se incluye una máscara *%year%* para que sea sustituida luego por cada año concreto.
- ***pathOutputDataFeeds***: ruta completa al directorio donde se van a escribir los ficheros con las características preseleccionadas. Se puede emplear la misma máscara para fijar el nombre en función del año concreto.
- ***campoAReemplazar***: es la máscara que se va a sustituir en el *pathDataFeeds* y *pathOutputDataFeeds* para reemplazarla por un año concreto. (En la implementación se ha utilizado *%year%*).

□ Clase **Division.java**

- ***pathContratosAnyos***: ruta completa al nombre de los ficheros de entrada con las características ya preseleccionadas. Se puede utilizar máscara para sustituir el año (típicamente este parámetro tendrá el mismo valor que *pathOutputDataFeeds* de Preselección, pues los ficheros de salida de esa fase se utilizan como entrada en esta).
- ***pathModificacionesContratos***: ruta completa de salida donde escribir los ficheros con las modificaciones de contratos. Se puede emplear máscara.

- ***pathNuevosContratos***: ruta completa de salida donde escribir los ficheros con información de nuevos contratos. Se puede emplear máscara.
- ***pathModificacionesContratosAgregados***: para estadísticas y elaboración de unos gráficos de la memoria se han agregado las modificaciones; en este parámetro se indica la ruta completa donde se escribe el fichero con la información de modificaciones de contratos agregadas para un año concreto. Se puede emplear máscara.
- ***pathNuevosContratosAgregados***: igual que el parámetro anterior pero aplicado a nuevos contratos.
- ***campoAReemplazar***: máscara que se emplea para identificar el año y sustituirla cuando aparece.

□ Clase **SeleccionConversion.java**

Los parámetros indicados aplican a las clases de implementación y experimentación dedicadas a la selección y conversión de los ficheros de contratos.

- ***pathFicherosCodificacion***: ruta (sólo el *path*) donde se van a escribir los ficheros que recogen la información de las características y los valores que se corresponde con cada literal codificado.
- ***pathTotales***: ruta (sólo el *path*) donde en la fase anterior de análisis se han generado unos ficheros de totales por cada característica empleando R. Los siguientes parámetros se especifica el nombre de cada fichero para algunas características concretas.
- ***ficheroTotalesAgencia***: nombre del fichero de texto con la información de agencias.
- ***ficheroTotalesClaimantProgramCode***: nombre del fichero de texto con la información de Claimant Program Code.
- ***ficheroTotalesContractionType***: nombre del fichero de texto con la información de Contraction Type.
- ***ficheroTotalesContractingOfficeAgencyId***: nombre del fichero de texto con la información de Contracting Office Agency Id.
- ***ficheroTotalesProductorServiceCode***: nombre del fichero de texto con la información de Productor Service Code.
- ***ficheroTotalesNaicsCode***: nombre del fichero de texto con la información de NAICS.

- ***ficheroTotalesReasonForModification***: nombre del fichero de texto con la información de Reason For Modification.
- ***pathFicheroNuevasOportunidades***: ruta completa de los ficheros de entrada con la información de las características preseleccionadas para nuevas contrataciones. Se puede emplear máscara. (Típicamente, los ficheros de entrada de esta clase son los de salida de nuevas contrataciones de la fase de División, por lo que será habitual que este parámetro tome el mismo valor que *pathNuevosContratos* en *Division.properties*).
- ***pathFicheroNuevasOportunidadesAgrupadas***: ruta completa de los ficheros de salida donde se van a escribir los ficheros resultantes con las características ya codificadas y agrupadas. Se puede utilizar máscara.
- ***campoAReemplazar***: máscara empleada para ser sustituida por el año.

□ Clase Recuento.java

Los parámetros aquí indicados aplican a las clases de implementación y experimentación dedicadas al recuentos de las tuplas de características y su aparición en la concesión de contratos.

- ***pathFicheroTuplas***: ruta completa del fichero de escritura donde se escribe por cada línea una tupla que aparece en las series de datos.
- ***pathFicheroTuplasContador***: ruta del fichero de salida donde se incluye una línea por cada tupla, que además de los campos de ésta incluye el número de observaciones distintas de 0 que aparecen para ella en el total de los años.
- ***pathFicheroNuevasOportunidadesAgrupadas***: ruta completa de los ficheros de entrada de nueva características codificadas y agrupadas. Se puede emplear máscara para ser sustituida por el año (estos ficheros de entrada son los generados en la fase de selección y conversión anterior por lo que este valor será similar al del parámetro *pathFicheroNuevasOportunidadesAgrupadas* del fichero de propiedades anterior).
- ***campoAReemplazar***: máscara que identifica al año para ser reemplazada.

□ Clase Completacion.java

Los parámetros aquí indicados aplican a las clases encargadas de completas las observaciones de las series temporales tanto en la parte de implementación como de experimentación.

- ***pathFicheroTuplas***: ruta completa al fichero con la colección de tuplas encontradas. Típicamente este fichero será el mismo que el generado en la fase de recuento anterior indicado por el parámetro *pathFicheroTuplas*.
- ***pathFicheroPlantillaAuxiliar***: ruta completa donde se va a generar el fichero auxiliar que sirva como plantilla de todas las observaciones que deberían aparecer de todas las tuplas durante un año.
- ***pathFicheroNuevasOportunidadesAgrupadas***: ruta completa de los ficheros de nuevas contrataciones codificados. Se emplea máscara para ser sustituida por el año. Estos ficheros son los generados en la fase de selección y conversión que se escriben en la ruta indicada por el parámetro *pathFicheroNuevasOportunidadesAgrupadas*.
- ***pathFicheroNuevasOportunidadesAgrupadasCompleto***: ruta completa donde se van a escribir los ficheros de salida definitivos, que incluyen las observaciones a 0 que faltaran. Se emplea máscara para identificar el año.
- ***campoAREemplazar***: máscara que parametriza el valor del año para ser sustituido.

□ Clase ConversionLSTM.java

- ***pathFicheroTuplas***: ruta completa al fichero con la colección de tuplas encontradas. Típicamente este fichero será el mismo que el generado en la fase de Recuento indicado por el parámetro *pathFicheroTuplas*.
- ***pathFicheroNuevasOportunidadesAgrupadasCompleto***: ruta de los ficheros de entrada con todas las observaciones completas. Se emplea máscara para el año. Se emplea máscara. Estos ficheros son los mismos que se generan en la fase de Completación y que se escriben en el ruta indicada por el parámetro *pathFicheroNuevasOportunidadesAgrupadasCompleto* de *Completacion.properties*.
- ***pathFicheroConversionLSTM***: ruta completa al único fichero de salida donde se incluye el valor de contratos asignados en cada uno de los 96 meses observados en la secuencia desde 2008 a 2015.
- ***campoAREemplazar***: máscara para poder ser sustituida por el año.

11. Bibliografía

□ Estado del Arte. *Big Data*.

[1] Marco Tulio Cicerón: “*La Adivinación; El Hado*”. (Editorial Folio, 2002. Traducción de la obra original *De Divinatione*, año 44 a.C).

[2] Cronología sobre el surgimiento del Big Data: <http://www.winshuttle.es/big-data-historia-cronologica> . (Último acceso en Junio de 2016).

[3] M. Cox, D. Ellsworth: “*Application-Controlled Demand Paging for Out-of-Core Visualization*”. NASA Ames Research Center, California, 1997. (Accesible en línea en <http://www.nas.nasa.gov/assets/pdf/techreports/1997/nas-97-010.pdf> . Último acceso en Junio de 2016).

[4] Cisco: “*Cisco Visual Networking Index: Forecast and Methodology, 2014–2019*”. Empresa Cisco, 2015. (Accesible en línea en http://www.cisco.com/c/en/us/solutions/collateral/service-provider/ip-ngn-ip-next-generation-network/white_paper_c11-481360.pdf . Último acceso en Junio de 2016).

[5] D. Laney: “*3D Data Management: Controlling Data Volume, Velocity, and Variety*”. Meta Group Inc., 2001. (Accesible en línea en <http://blogs.gartner.com/doug-laney/files/2012/01/ad949-3D-Data-Management-Controlling-Data-Volume-Velocity-and-Variety.pdf> . Último acceso en Junio de 2016).

[6] Sitio web oficial de Internet Live Stats: <http://www.internetlivestats.com/one-second> . (Último acceso en Junio de 2016).

[7] M. A. u. d. Khan, M. F. Uddin and N. Gupta: “*Seven V's of Big Data understanding Big Data to extract value*”. American Society for Engineering Education, Connecticut. (Accesible en línea en <http://www.asee.org/documents/zones/zone1/2014/Professional/PDFs/113.pdf> . Último acceso en Junio de 2016).

[8] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, I. Stoica: “*Spark: Cluster Computing with Working Sets*”. University of California, Berkeley. 2010. (Accesible en línea en <https://amplab.cs.berkeley.edu/wp-content/uploads/2011/06/Spark-Cluster-Computing-with-Working-Sets.pdf> . Último acceso en Junio de 2016).

[9] Sitio web oficial de las publicaciones del Proyecto de investigación Stratosphere: <http://stratosphere.eu/project/publications> . (Último acceso en Junio de 2016).

[10] R. Xin, J. Rosen: “*Project Tungsten: Bringing Apache Spark Closer to Bare Metal*”. Empresa Databricks, 2015. (Accesible en línea en <https://databricks.com/blog/2015/04/28/project-tungsten-bringing-spark-closer-to-bare-metal.html> . Último acceso en Junio de 2016).

[11] F. Hüske: “*Juggling with Bits and Bytes*”. Apache Flink, 2015. (Accesible en línea en <http://flink.apache.org/news/2015/05/11/Juggling-with-Bits-and-Bytes.html> . Último acceso en Junio de 2016).

[12] T. Vasiloudis: “*FlinkML: Vision and Roadmap*”, 2015. Accesible en línea en: <https://cwiki.apache.org/confluence/display/FLINK/FlinkML%3A+Vision+and+Roadmap> . (Último acceso en Junio de 2016).

[13] J. Dean, S. Ghemawat: “*MapReduce: simplified data processing on large clusters*”. Google Inc., 2008. (Accesible en línea en: <http://static.googleusercontent.com/media/research.google.com/es//archive/mapreduce-osdi04.pdf> . Último acceso en Junio de 2016).

[14] K. Leboeuf: “*2016 Update: What happens in one Internet minute?*”. Excelacom Inc., 2016. (Accesible en línea en: <http://www.excelacom.com/resources/blog/2016-update-what-happens-in-one-internet-minute> . Último acceso en Junio de 2016).

[15] Institute of Electrical and Electronics Engineers: “16326-2009 - ISO/IEC/IEEE Systems and Software Engineering--Life Cycle Processes--Project Management”. 2009.

[16] Licencia de Java: <http://www.oracle.com/technetwork/java/javase/terms/license/index.html> (Último acceso en Junio de 2016).

[17] M. Turck, J. Hao: “*Big Data Landscape 2016 3rd Version*”. FirstMark Capital, 2016. Accesible en línea en: <http://mattturck.com/wp-content/uploads/2016/03/Big-Data-Landscape-2016-v18-FINAL.png> . Último acceso en Junio de 2016).

[18] Gartner, Inc: “*Gartner's 2015 Hype Cycle for Emerging Technologies Identifies the Computing Innovations That Organizations Should Monitor*”. Gartner, Connecticut, 2015. (Accesible en línea en <http://www.gartner.com/newsroom/id/3114217> . Último acceso en Junio de 2016).

[19] Gartner, Inc: “*Gartner's 2014 Hype Cycle for Emerging Technologies Maps the Journey to Digital Business*”. Gartner, Connecticut, 2014. (Accesible en línea en <http://www.gartner.com/newsroom/id/2819918> . Último acceso en Junio de 2016).

□ Contratación Pública

[20] Sitio web oficial de la iniciativa Usa Spending: <https://www.usaspending.gov> . (Último acceso en Junio de 2016).

[21] Sitio web oficial de la iniciativa Public Spending: <http://publicspending.net> . (Último acceso en Junio de 2016).

[22] J.M. Álvarez Rodríguez: “*Métodos semánticos de reutilización de datos abiertos enlazados en las licitaciones públicas*”. Tesis Doctoral, Universidad de Oviedo, 2012.

[23] Usaspending.gov: “Data Downloads Data Dictionary”. 2015. (Accesible en línea: <https://www.usaspending.gov/about/PublishingImages/Pages/TheData/USAspending.gov%20Data%20Dictionary.pdf> . Último acceso en Junio de 2016).

[24] Sitio web oficial del North American Industry Classification System: <http://www.census.gov/eos/www/naics/> . (Último acceso en Junio de 2016).

□ Herramientas

[25] Sitio web oficial de Apache Flink: <https://flink.apache.org> . (Último acceso en Junio de 2016).

[26] Sitio web oficial de Apache Hadoop: <http://hadoop.apache.org> . (Último acceso en Junio de 2016).

[27] Sitio web oficial de Apache Mahout: <http://mahout.apache.org> . (Último acceso en Junio de 2016).

[28] Sitio web oficial de MongoDB en español: <https://www.mongodb.com/es> . (Último acceso en Junio de 2016).

[29] Sitio web oficial de Apache Spark: <http://spark.apache.org> . (Último acceso en Junio de 2016).

[30] Sitio web oficial de TensorFlow: <https://www.tensorflow.org> . (Último acceso en Junio de 2016).

[31] Sitio web oficial de Keras: <http://keras.io>. (Último acceso en Junio de 2016).

[32] Sitio web oficial de R: <https://www.r-project.org> . (Último acceso en Junio de 2016).

[33] Sitio web oficial de R Studio: <https://www.rstudio.com> . (Último acceso en Junio de 2016).

[34] Sitio web oficial de Java/R Interface: <https://rforge.net/JRI> . (Último acceso en Junio de 2016).

[35] Sitio web oficial de Docker: <https://www.docker.com> . (Último acceso en Junio de 2016).

[36] Sitio web oficial de Cloudera: <http://www.cloudera.com>. (Último acceso en Junio de 2016).

[37] Sitio web oficial de Google Cloud Platform: <https://cloud.google.com> . (Último acceso en Junio de 2016).

[38] Sitio web oficial de Hortonworks: <http://hortonworks.com> . (Último acceso en Junio de 2016).

- [39] Sitio web oficial de Apache Kafka: <http://kafka.apache.org> . (Último acceso en Junio de 2016).
- [40] Sitio web oficial de Apache Storm: <http://storm.apache.org> . (Último acceso en Junio de 2016).
- [41] Transformaciones implementadas en Apache Flink: https://ci.apache.org/projects/flink/flink-docs-release-0.7/dataset_transformations.html . (Último acceso en Junio de 2016).
- [42] Sitio web oficial de Python: <https://www.python.org> . (Último acceso en Junio de 2016).
- [43] Sitio web oficial de Java es español: <http://www.java.com/es/> . (Último acceso en Junio de 2016).
- [44] Sitio web oficial de Eclipse Mars: <https://eclipse.org/mars> . (Último acceso en Junio de 2016).
- [45] FlinkCEP: <https://ci.apache.org/projects/flink/flink-docs-master/apis/streaming/libs/cep.html> . (Último acceso en Junio de 2016).
- [46] Gelly: <https://ci.apache.org/projects/flink/flink-docs-master/apis/batch/libs/gelly.html> . (Último acceso en Junio de 2016).
- [47] Table API and SQL: <https://ci.apache.org/projects/flink/flink-docs-master/apis/table.html> . (Último acceso en Junio de 2016).
- [48] Documentación de la clase QuantileDiscretizer de Apache Spark. Accesible en línea en: <https://spark.apache.org/docs/1.6.1/api/java/org/apache/spark/ml/feature/QuantileDiscretizer.html> . (Último acceso en Junio de 2016).
- [49] Documentación de la clase Bucketizer de Apache Spark. Accesible en línea en: <https://spark.apache.org/docs/1.4.1/api/java/org/apache/spark/ml/feature/Bucketizer.html> . (Último acceso en Junio de 2016).
- [50] Documentación de la clase Normalizer de Apache Spark. Accesible en línea en: <https://spark.apache.org/docs/1.4.0/api/java/org/apache/spark/mllib/feature/Normalizer.html> . (Último acceso en Junio de 2016).
- [51] “Feature Extraction and Transformation with spark.mllib”: <http://spark.apache.org/docs/latest/mllib-feature-extraction.html> . (Último acceso en Junio de 2016).
- [52] Documentación de la clase StandardScaler de Apache Spark. Accesible en línea en: <https://spark.apache.org/docs/1.4.0/api/java/index.html?org/apache/spark/mllib/feature/StandardScaler.html>

□ Redes Neuronas

[53] C. Olah: “Understanding LSTM Networks”. 2015. (Accesible en línea en: <http://colah.github.io/posts/2015-08-Understanding-LSTMs> . Último acceso en Junio de 2016).

[54] W. S. McCulloch , W. Pitts: “*A Logical Calculus Of The Ideas Immanent In Nervous Activity*”. Bulletin of Mathematical Biophysics, Vol. 5, pp. 115-133, 1943. (Versión accesible en línea en: <http://deeplearning.cs.cmu.edu/pdfs/McCulloch.and.Pitts.pdf> . Último acceso en Junio de 2016).

[55] S. Hochreiter, J. Schmidhuber: “*Long Short-Term Memory*”. Neural Computation 9(8). 1997. (Versión accesible en línea en: http://deeplearning.cs.cmu.edu/pdfs/Hochreiter97_lstm.pdf . Último acceso en Junio de 2016).

[56] F. A. Gers, J. Schmidhuber , F. Cummins: “*Learning to Forget: Continual Prediction with LSTM*”. Technical Report IDSIA-01-99, 1999.

[57] Tutorial en Deeplearning.net sobre redes LSTM. Accesible en línea en: <http://deeplearning.net/tutorial/lstm.html> . (Último acceso en Junio de 2016).

[58] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov: “*Dropout: A Simple Way to Prevent Neural Networks from Overfitting*”. Journal of Machine Learning Research 15, 2014. (Accesible en línea en: <http://www.cs.toronto.edu/~rsalakhu/papers/srivastava14a.pdf> . Último acceso en Junio de 2016).

[59] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mane, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viegas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, X. Zheng: “*TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems*”. Google Research, 2015. (Accesible en línea en <http://download.tensorflow.org/paper/whitepaper2015.pdf> . Último acceso en Junio de 2016).

[60] P. Isasi, I. M. Galván: “*Redes de Neuronas Artificiales. Un enfoque Práctico*”. Pearson Prentice Hall, 2004.

□ Regresión Lineal

[61] D. M. Diez, C. D. Barr, M. Cetinkaya-Rundel : “*OpenIntro Statistics*”. OpenIntro, 2016, 3ª Edición.

[62] Función lm del paquete stats de R. Accesible documentación en línea en: <https://stat.ethz.ch/R-manual/R-devel/library/stats/html/lm.html> (Último acceso en Junio de 2016).

[63] Función cor del paquete stats de R. Accesible documentación en línea en: <https://stat.ethz.ch/R-manual/R-devel/library/stats/html/cor.html> . (Último acceso en Junio de 2016).

[64] X. Yan, X. Gang Su: “*Linear Regression Analysis: Theory and Computing*”. World Scientific, 2009.

□ Series temporales y Modelos ARIMA

[65] T. Villagarcía: “*Apuntes sobre series temporales*”. (Accesible en línea en http://www.est.uc3m.es/esp/nueva_docencia/leganes/ing_industrial/estadistica_industria/doc_grupo1/archivos/Apuntes%20de%20series.pdf . Último acceso en Junio de 2016).

[66] D. Peña: “*Análisis de series temporales*”. (Alianza Editorial, Madrid, 2010).

[67] Paquete Forecast de R: Función auto.arima (Creado por R. J. Hyndman). Accesible la documentación de la función en línea en: <http://www.inside-r.org/packages/cran/forecast/docs/auto.arima> (Último acceso en Junio de 2016).

[68] R. Hyndman, Y. Khandakar: “*Automatic Time Series Forecasting: The forecast Package for R*”. Journal of Statistical Software, 2008.

[69] F. Daza, N. Ceular, J. M. Caridad y Ocerin: “*Desarrollo de una metodología alternativa a partir del modelo ARIMA y redes neuronales para la predicción del consumo de agua en zonas urbanas*”. Universidad de Córdoba, 2008. (Accesible en línea en: <http://www.iiis.org/cds2008/cd2008csc/cisci2008/paperspdf/c407yt.pdf> . Último acceso en Junio de 2016).

[70] D. Quintana, R. Gimeno, P. Isasi: “Detección de inercia sectorial en salidas a bolsa mediante modelos ARIMA y redes neuronales”. Economía y Administración N° 65, pp. 29 – 53. 2005.